

# Deep Learning

An Introduction

Mirko Bunse, Quentin Fühling, and Vukan Jevtic

Grad School on Astro-Particle Physics (Jan 18<sup>th</sup>–23<sup>rd</sup>, 2026)

Partner institutions:



Institutionally funded by:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

Ministerium für  
Kultur und Wissenschaft  
des Landes Nordrhein-Westfalen



## Further Reading

**Goodfellow, Bengio, and Courville, *Deep Learning*, 2016:**

- ▶ principled and rigorous approach
- ▶ great technical coverage





## Further Reading

### **Goodfellow, Bengio, and Courville, *Deep Learning*, 2016:**

- ▶ principled and rigorous approach
- ▶ great technical coverage

### **Bishop and Bishop, *Deep Learning - Foundations and Concepts*, 2024:**

- ▶ even more basics
- ▶ more advanced topics, like transformers



## Further Reading

### **Goodfellow, Bengio, and Courville, *Deep Learning*, 2016:**

- ▶ principled and rigorous approach
- ▶ great technical coverage

### **Bishop and Bishop, *Deep Learning - Foundations and Concepts*, 2024:**

- ▶ even more basics
- ▶ more advanced topics, like transformers

### **Erdmann et al., *Deep Learning for Physics Research*, 2021:**

- ▶ physics-oriented examples and exercises
- ▶ (some) coverage of uncertainties and custom loss functions



## Further Reading

### **Goodfellow, Bengio, and Courville, *Deep Learning*, 2016:**

- ▶ principled and rigorous approach
- ▶ great technical coverage

### **Bishop and Bishop, *Deep Learning - Foundations and Concepts*, 2024:**

- ▶ even more basics
- ▶ more advanced topics, like transformers

### **Erdmann et al., *Deep Learning for Physics Research*, 2021:**

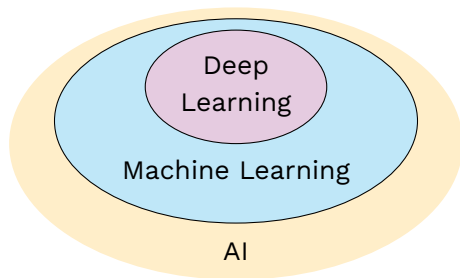
- ▶ physics-oriented examples and exercises
- ▶ (some) coverage of uncertainties and custom loss functions

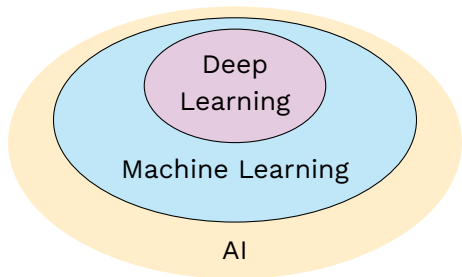
### **Lippe, *UvA Deep Learning Tutorials*, 2023:**

- ▶ <https://uvadlc-notebooks.readthedocs.io>



# Introduction / Machine Learning





**Machine learning** = data ○ model ○ fit



# Supervised Learning



**Target:** any quantity  $Y$  we want to predict (costly or impossible to measure)

**Feature:** any quantity  $X_i$  we compute from observable quantities

**Training Data:**  $D = \{ (x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : 1 \leq i \leq m \}$

# Supervised Learning



**Target:** any quantity  $Y$  we want to predict (costly or impossible to measure)

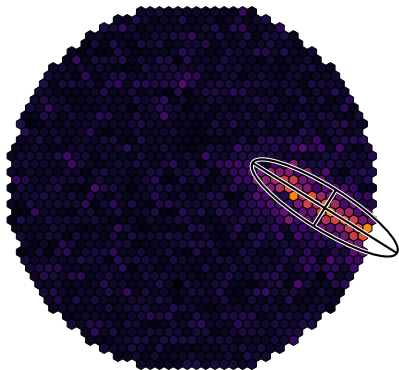
**Feature:** any quantity  $X_i$  we compute from observable quantities

**Training Data:**  $D = \{ (x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : 1 \leq i \leq m \}$

target	feature	feature	
$Y$	$X_1$	$X_2$	
+1	1.3	A	...
-1	-0.2	B	
+1	0.8	A	
	⋮		⋱

## Structured data:

- ▶ tabular representation
- ▶  $X_i$  facilitate the prediction of  $Y$ ,  
e.g., through well-designed preprocessing



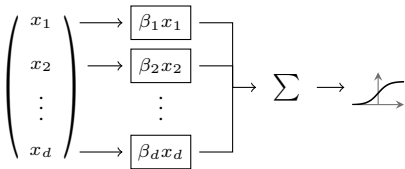
```
df = fact.io.read_data( # pandas.DataFrame
    "gamma_simulations_facttools_dl2.hdf5",
    key = "events"
)

X = df[[ # select features
    "length", # -> shape (n_events, n_features)
    "width",
    "num_islands",
    "num_pixel_in_shower",
    # ...
]].to_numpy()

y = df["corsika_event_header_total_energy"]

clf = sklearn.ensemble.RandomForestClassifier()
clf.fit(X, y)
```

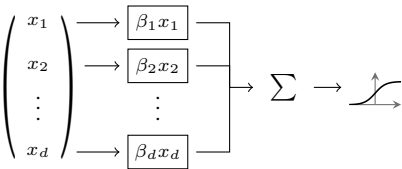
# Structured Data



**Logistic Regression:**

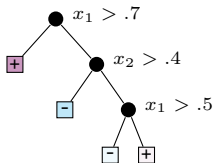
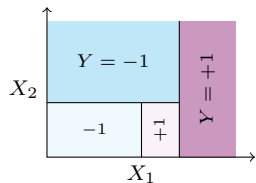
$$\hat{\mathbb{P}}_{\beta}(Y = +1 \mid X = x) = \frac{e^{\langle \beta, x \rangle}}{1 + e^{\langle \beta, x \rangle}}$$

# Structured Data



## Logistic Regression:

$$\hat{\mathbb{P}}_{\beta}(Y = +1 \mid X = x) = \frac{e^{\langle \beta, x \rangle}}{1 + e^{\langle \beta, x \rangle}}$$

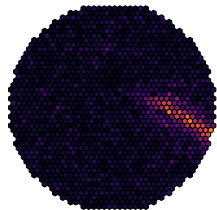


## Decision Trees:

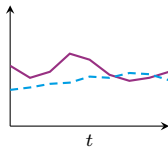
- ▶ recursively split  $\mathcal{X}$
- ▶ boost performance through ensembling

These models perform very well 🎉 (if structure permits 😞)

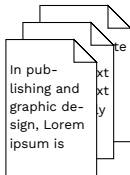
# Unstructured Data



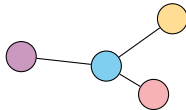
**images**



**time series**

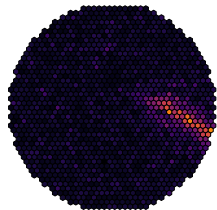


**texts**

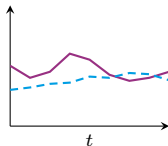


**graphs**

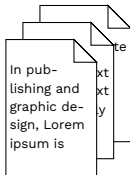
# Unstructured Data



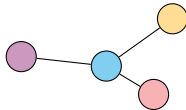
images



time series



texts

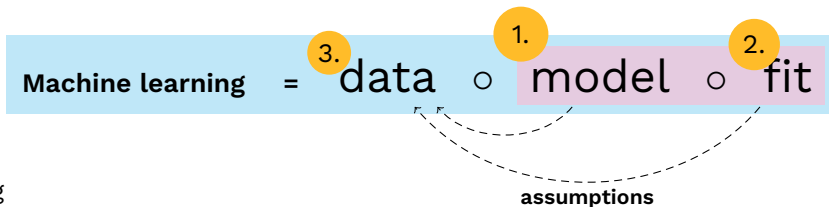


graphs

**Deep Learning** learns features as a part of the model 🚀

- ▶ no manual feature-engineering necessary 🙌
- ▶ instead, architecture optimization and more data are needed 😞

# Agenda



1. Modeling
  2. Fitting
  3. Data and Assumptions
  4. Concluding Remarks
- + Hands-On Exercises (Tue ~ 45 min, Thu ~ 90 min)





# Modeling



**Polynomial Regression:**  $y = f_{\beta}(x) + \epsilon$ , where  $f_{\beta}(x) = \sum_{i=0}^n \langle \beta_i, x^i \rangle$

- ▶  $y \in \mathbb{R}$ ,  $x \in \mathbb{R}^d$ , and  $\beta_i \in \mathbb{R}^d$
- ▶  $\langle a, b \rangle = \sum_{j=1}^d a_j \cdot b_j$  is the scalar product



**Polynomial Regression:**  $y = f_{\beta}(x) + \epsilon$ , where  $f_{\beta}(x) = \sum_{i=0}^n \langle \beta_i, x^i \rangle$

- ▶  $y \in \mathbb{R}$ ,  $x \in \mathbb{R}^d$ , and  $\beta_i \in \mathbb{R}^d$
- ▶  $\langle a, b \rangle = \sum_{j=1}^d a_j \cdot b_j$  is the scalar product
- ▶ typical loss:  $\mathcal{L}_D(\beta) = \sum_{i=1}^m (y_i - f_{\beta}(x_i))^2$  where  $(x_i, y_i) \in D$  and  $D = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : 1 \leq i \leq m\}$  is the training set



**Logistic Regression:**  $\hat{y} = \arg \max_{i \in \{1, 2, \dots, C\}} \hat{\mathbb{P}}_{\beta}(Y = i \mid X = x)$



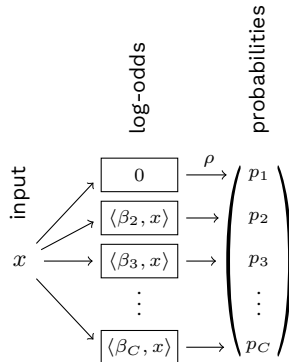
**Logistic Regression:**  $\hat{y} = \arg \max_{i \in \{1, 2, \dots, C\}} \underbrace{\hat{\mathbb{P}}_{\beta}(Y = i \mid X = x)}_{= \rho(\langle \beta_i, x \rangle)}$

$$\text{where } \rho(v_i) = \begin{cases} \frac{1}{1 + \sum_{j=2}^k e^{v_j}} & i = 1 \\ \frac{e^{v_i}}{1 + \sum_{j=2}^k e^{v_j}} & i \in \{2, 3, \dots, C\} \end{cases}$$



**Logistic Regression:**  $\hat{y} = \arg \max_{i \in \{1, 2, \dots, C\}} \underbrace{\hat{\mathbb{P}}_{\beta}(Y = i \mid X = x)}_{= \rho(\langle \beta_i, x \rangle)}$

$$\text{where } \rho(v_i) = \begin{cases} \frac{1}{1 + \sum_{j=2}^k e^{v_j}} & i = 1 \\ \frac{e^{v_i}}{1 + \sum_{j=2}^k e^{v_j}} & i \in \{2, 3, \dots, C\} \end{cases}$$



The **soft-max** operation  $\rho$  projects to the unit simplex  $\{p \in \mathbb{R}^C : p_i \geq 0, 1 = \sum_{i=1}^C p_i\}$



**Motivation:** the Logistic Regression represents **linear models** of the **log-odds**.

$$\log \frac{\mathbb{P}(Y = 2 \mid X = x)}{\mathbb{P}(Y = 1 \mid X = x)} = \langle \beta_2, x \rangle + \epsilon \stackrel{?}{>} 0$$

$$\log \frac{\mathbb{P}(Y = 3 \mid X = x)}{\mathbb{P}(Y = 1 \mid X = x)} = \langle \beta_3, x \rangle + \epsilon \stackrel{?}{>} 0$$

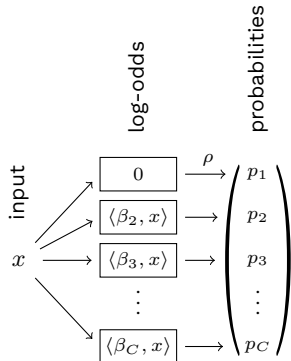
...

$$\log \frac{\mathbb{P}(Y = k \mid X = x)}{\mathbb{P}(Y = 1 \mid X = x)} = \langle \beta_C, x \rangle + \epsilon \stackrel{?}{>} 0$$



## Synopsis:

- ▶ **Polynomial Regression** =  
a linear model of exponentiated inputs  $x^i$
- ▶ **Logistic Regression** =  
a linear model of the log-odds
- ▶ The **soft-max** operation maps these log-odds  
to (estimates of) class probabilities



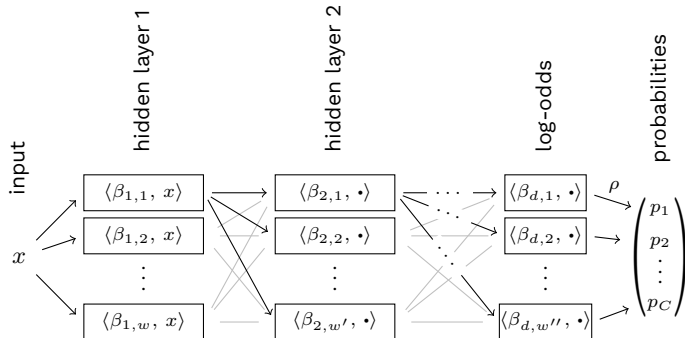


# Deep Networks



**Deep Nets:** use multiple (logistic regression-like) layers

- ▶ learnable linear combinations  $\langle \beta, \cdot \rangle$

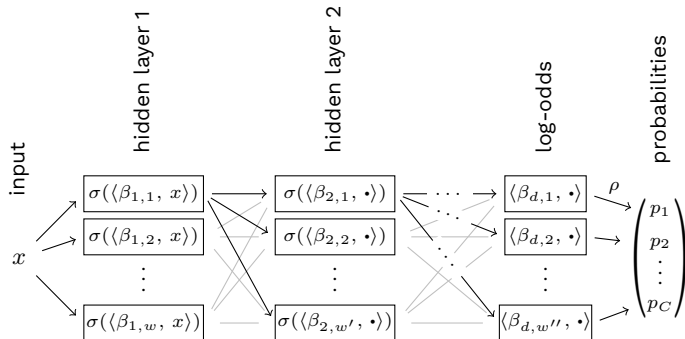
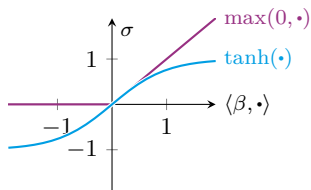


# Deep Networks



**Deep Nets:** use multiple (logistic regression-like) layers

- ▶ learnable linear combinations  $\langle \beta, \cdot \rangle$
- ▶ non-linear activations  $\sigma$



# Universal Approximation



**Density:** A family  $G$  of models can approximate any function  $f \in C(\mathbb{R}^n)$ ,  
if  $\forall \varepsilon > 0$ , compact  $K \subseteq \mathbb{R}^n$ ,  $\exists g \in G$ , such that

$$\max_{x \in K} \|f(x) - g(x)\| < \varepsilon$$

<sup>1</sup> Pinkus, “Approximation theory of the MLP model in neural networks”, 1999.

<sup>2</sup> Kidger and Lyons, “Universal approximation with deep narrow networks”, 2020.

# Universal Approximation



**Density:** A family  $G$  of models can approximate any function  $f \in C(\mathbb{R}^n)$ , if  $\forall \varepsilon > 0$ , compact  $K \subseteq \mathbb{R}^n$ ,  $\exists g \in G$ , such that

$$\max_{x \in K} \|f(x) - g(x)\| < \varepsilon$$

- ▶ **One hidden layer of arbitrary width** is dense iff  $\sigma$  is non-polynomial.<sup>1</sup>
- ▶ **Arbitrarily deep nets with minimum width**  $d + C + 2$  are dense.<sup>2</sup>

<sup>1</sup> Pinkus, “Approximation theory of the MLP model in neural networks”, 1999.

<sup>2</sup> Kidger and Lyons, “Universal approximation with deep narrow networks”, 2020.

# Universal Approximation



**Density:** A family  $G$  of models can approximate any function  $f \in C(\mathbb{R}^n)$ , if  $\forall \varepsilon > 0$ , compact  $K \subseteq \mathbb{R}^n$ ,  $\exists g \in G$ , such that

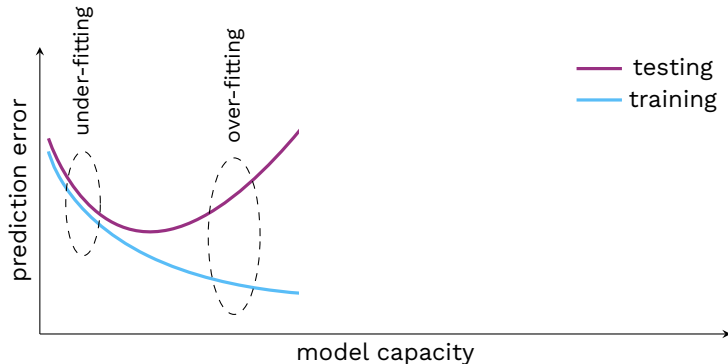
$$\max_{x \in K} \|f(x) - g(x)\| < \varepsilon$$

- ▶ **One hidden layer of arbitrary width** is dense iff  $\sigma$  is non-polynomial.<sup>1</sup>
- ▶ **Arbitrarily deep nets with minimum width**  $d + C + 2$  are dense.<sup>2</sup>
- ▶ Deep nets are often more *efficient* approximators than wide shallow nets.
- ▶ Density does not imply the existence of a learning algorithm to select  $g$  from  $G$

<sup>1</sup> Pinkus, “Approximation theory of the MLP model in neural networks”, 1999.

<sup>2</sup> Kidger and Lyons, “Universal approximation with deep narrow networks”, 2020.

# Over- and Underfitting



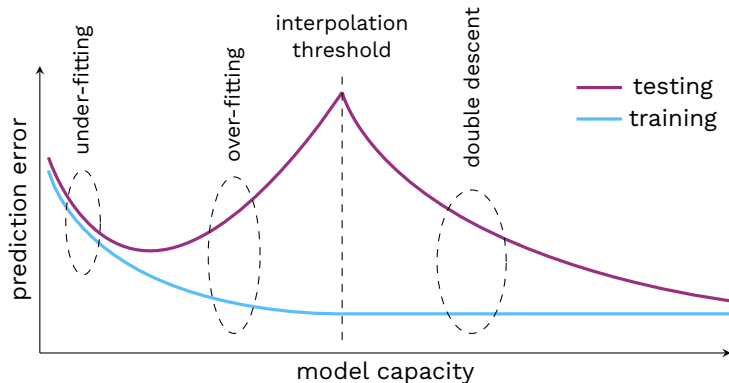
## Under-Fitting:

- ▶ approximation
- ▶ high bias, low variance

## Over-Fitting:

- ▶ memorization
- ▶ low bias, high variance

# Over- and Underfitting



## Under-Fitting:

- ▶ approximation
- ▶ high bias, low variance

## Over-Fitting:

- ▶ memorization
- ▶ low bias, high variance

## Double Descent:

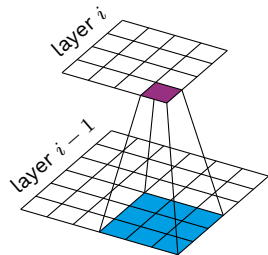
- ▶ interpolation<sup>3</sup>

<sup>3</sup> Belkin et al., “Reconciling modern machine-learning practice and the classical bias-variance trade-off”, 2019

# Inductive Biases



**Convolution:**  $S(i, j) = (K * I)(i, j) = \sum_{m, n} I(i - m, j - n) \cdot K(m, n)$

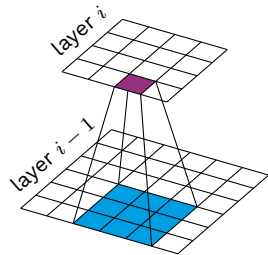




# Inductive Biases



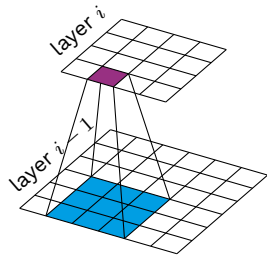
**Convolution:**  $S(i, j) = (K * I)(i, j) = \sum_{m, n} I(i - m, j - n) \cdot K(m, n)$



# Inductive Biases



**Convolution:**  $S(i, j) = (K * I)(i, j) = \sum_{m, n} I(i - m, j - n) \cdot K(m, n)$

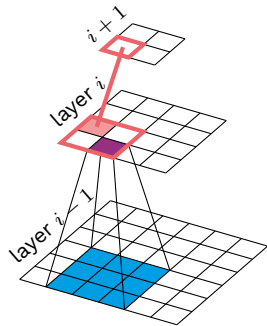


# Inductive Biases



**Convolution:**  $S(i, j) = (K * I)(i, j) = \sum_{m, n} I(i - m, j - n) \cdot K(m, n)$

**Pooling:** only maintain the **maximum** of each neighborhood.



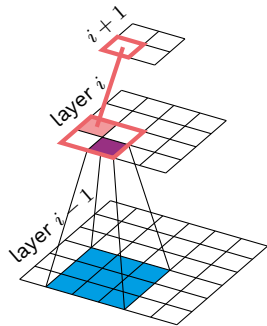
# Inductive Biases



**Convolution:**  $S(i, j) = (K * I)(i, j) = \sum_{m, n} I(i - m, j - n) \cdot K(m, n)$

**Pooling:** only maintain the **maximum** of each neighborhood.

- ▶ translation invariance
- ▶ sparse interactions
- ▶ parameter sharing

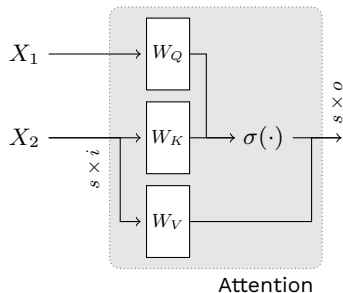


In general, specialized layers are used to introduce **biases** that suit the data.

# Inductive Biases



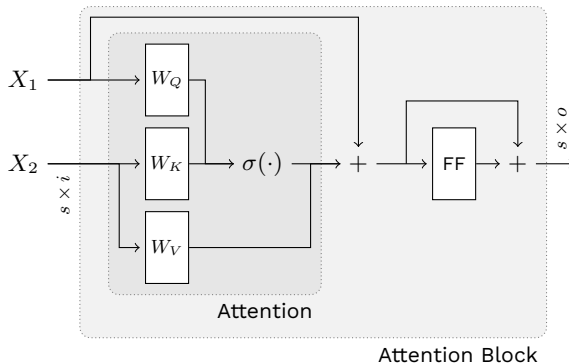
$\text{Attention}(X_1, X_2) = \sigma((X_1 W_Q)(X_2 W_K)^\top) X_2 W_V$  explicitly models interactions.



# Inductive Biases



$\text{Attention}(X_1, X_2) = \sigma((X_1 W_Q)(X_2 W_K)^\top) X_2 W_V$  explicitly models interactions.



# Modeling



## Synopsis:

- ▶ **Deep Nets** use layers of increasingly abstract representations
- ▶ **Layers** consist of linear parameters and non-linear activations

# Modeling



## Synopsis:

- ▶ **Deep Nets** use layers of increasingly abstract representations
- ▶ **Layers** consist of linear parameters and non-linear activations
- ▶ **Model Capacity** should consider sample sizes (over-/under-fitting)
- ▶ **Inductive Biases** facilitate learning



# Modeling



## Synopsis:

- ▶ **Deep Nets** use layers of increasingly abstract representations
- ▶ **Layers** consist of linear parameters and non-linear activations
- ▶ **Model Capacity** should consider sample sizes (over-/under-fitting)
- ▶ **Inductive Biases** facilitate learning

## Practical Recommendations:

- ▶ **Build on Existing Solutions** for similar problems

# Modeling



## Synopsis:

- ▶ **Deep Nets** use layers of increasingly abstract representations
- ▶ **Layers** consist of linear parameters and non-linear activations
- ▶ **Model Capacity** should consider sample sizes (over-/under-fitting)
- ▶ **Inductive Biases** facilitate learning

## Practical Recommendations:

- ▶ **Build on Existing Solutions** for similar problems
- ▶ **Extensively Tune** the hyper-parameters (# layers, # features per layer, ...)

# Modeling



## Synopsis:

- ▶ **Deep Nets** use layers of increasingly abstract representations
- ▶ **Layers** consist of linear parameters and non-linear activations
- ▶ **Model Capacity** should consider sample sizes (over-/under-fitting)
- ▶ **Inductive Biases** facilitate learning

## Practical Recommendations:

- ▶ **Build on Existing Solutions** for similar problems
- ▶ **Extensively Tune** the hyper-parameters (# layers, # features per layer, ...)
- ▶ **Assumptions > Depth** hence, prioritize baseline methods



# Fitting

# Empirical Risk Minimization



## Notation:

- ▶  $h_\beta : \mathcal{X} \rightarrow \mathbb{R}^C$  is our model, parametrized by  $\beta \in \mathbb{R}^B$  (fixed architecture)
- ▶  $\ell(h_\beta(x), y)$  measures the deviation between  $h_\beta(x)$  and  $y$

# Empirical Risk Minimization



## Notation:

- ▶  $h_\beta : \mathcal{X} \rightarrow \mathbb{R}^C$  is our model, parametrized by  $\beta \in \mathbb{R}^B$  (fixed architecture)
- ▶  $\ell(h_\beta(x), y)$  measures the deviation between  $h_\beta(x)$  and  $y$

**Ultimate Goal:** minimize the *expected* risk:

$$R(h_\beta, \ell) = \mathbb{E}_{(x,y) \sim \mathbb{P}}(\ell(h_\beta(x), y)) = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{P}(X = x, Y = y) \cdot \ell(h_\beta(x), y) \, dx \, dy$$

# Empirical Risk Minimization



## Notation:

- ▶  $h_\beta : \mathcal{X} \rightarrow \mathbb{R}^C$  is our model, parametrized by  $\beta \in \mathbb{R}^B$  (fixed architecture)
- ▶  $\ell(h_\beta(x), y)$  measures the deviation between  $h_\beta(x)$  and  $y$

**Ultimate Goal:** minimize the *expected* risk:

$$R(h_\beta, \ell) = \mathbb{E}_{(x,y) \sim \mathbb{P}}(\ell(h_\beta(x), y)) = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{P}(X = x, Y = y) \cdot \ell(h_\beta(x), y) \, dx \, dy$$

**Approach:** approximate  $R(h_\beta, \ell)$  empirically with the training data  $D$ :

$$\widehat{R}_D(h_\beta, \ell) = \frac{1}{m} \sum_{i=1}^m \ell(h_\beta(x_i), y_i) \xrightarrow{m \rightarrow \infty} R(h_\beta, \ell)$$

and choose  $\beta^* = \arg \min_{\beta \in \mathbb{R}^B} \widehat{R}_D(h_\beta, \ell)$ .



**Mean Squared Error:**  $\ell(h(x), y) = \|h(x) - y\|_2^2$

**Cross Entropy / Logistic Loss:**  $\ell'(h(x), y) = - \sum_{i=1}^C \delta_{y=i} \log([h(x)]_i)$





**Mean Squared Error:**  $\ell(h(x), y) = \|h(x) - y\|_2^2$

**Cross Entropy / Logistic Loss:**  $\ell'(h(x), y) = - \sum_{i=1}^C \delta_{y=i} \log([h(x)]_i)$

**Proper Scoring Rule:** any  $\ell : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$  for which  $\arg \min_{h \in \mathcal{H}} R(h; \ell) = \mathbb{P}(Y | X)$ .

- ▶ cross entropy is proven to be such a loss function
- ▶ hence, ERM with cross entropy readily learns  $\mathbb{P}(Y | X)$  🚀

# Empirical Risk Minimization (Revisited)



**Ultimate Goal:** minimize the *expected* risk:

$$R(h_\beta, \ell) = \mathbb{E}_{(x,y) \sim \mathbb{P}}(\ell(h_\beta(x), y))$$

**ERM:** approximate  $R(h_\beta, \ell)$  empirically with the training data  $D$ :

$$\hat{R}_D(h_\beta, \ell) = \frac{1}{m} \sum_{i=1}^m \ell(h_\beta(x_i), y_i) \xrightarrow{m \rightarrow \infty} R(h_\beta, \ell)$$

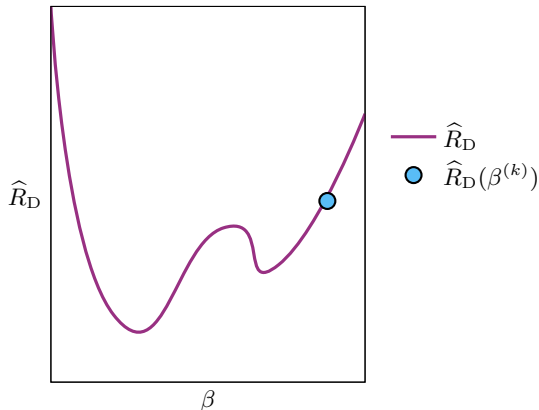
and choose  $\beta^* = \arg \min_{\beta \in \mathbb{R}^B} \hat{R}_D(h_\beta, \ell)$ .

# Stochastic First-Order Optimization



## Ideas:

- $\hat{R}_D(h_\beta, \ell)$  is just a function to be minimized

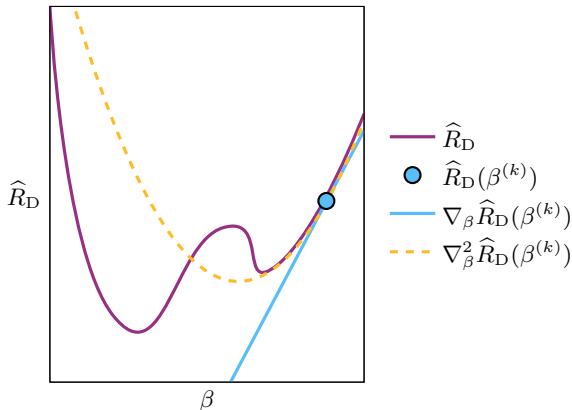


# Stochastic First-Order Optimization



## Ideas:

- ▶  $\hat{R}_D(h_\beta, \ell)$  is just a function to be minimized
- ▶ use gradient information to reduce  $\hat{R}_D(h_\beta, \ell)$  until  $\beta^*$  is found.
- ▶ ignore higher-order derivatives to save computation time.

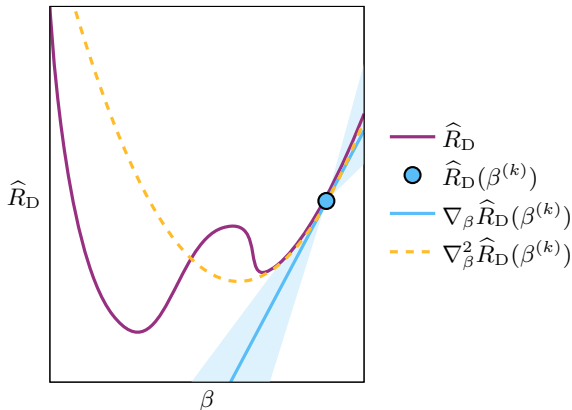


# Stochastic First-Order Optimization



## Ideas:

- ▶  $\hat{R}_D(h_\beta, \ell)$  is just a function to be minimized
- ▶ use gradient information to reduce  $\hat{R}_D(h_\beta, \ell)$  until  $\beta^*$  is found.
- ▶ ignore higher-order derivatives to save computation time.
- ▶ introduce randomness into the gradients to improve convergence.



# Stochastic First-Order Optimization



**Stochastic Gradient Descent (SGD):** in each step  $k$ , reduce the risk  $\widehat{R}_D(h_\beta, \ell)$  w.r.t. a *single, random* example.

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \alpha^{(k)} \nabla_{\beta} \ell \left( h(x_{i^{(k)}}), y_{i^{(k)}} \right) \text{ where } \begin{cases} \beta^{(k)} & \text{the parameter vector of } h \\ \alpha^{(k)} & \text{the step size} \\ (x_{i^{(k)}}), y_{i^{(k)}} & \text{the example} \end{cases}$$

# Stochastic First-Order Optimization



**Stochastic Gradient Descent (SGD):** in each step  $k$ , reduce the risk  $\widehat{R}_D(h_\beta, \ell)$  w.r.t. a *single, random* example.

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \alpha^{(k)} \nabla_\beta \ell \left( h(x_{i^{(k)}}), y_{i^{(k)}} \right) \quad \text{where} \quad \begin{cases} \beta^{(k)} & \text{the parameter vector of } h \\ \alpha^{(k)} & \text{the step size} \\ (x_{i^{(k)}}), y_{i^{(k)}} & \text{the example} \end{cases}$$

**Full Gradient Descent (GD):** in each step  $k$ , reduce  $\widehat{R}_D(h_\beta, \ell)$  w.r.t. *all* examples.

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \alpha^{(k)} \nabla_\beta \widehat{R}_D(h_\beta, \ell) = \beta^{(k)} - \alpha^{(k)} \frac{1}{m} \sum_{i=1}^m \nabla_\beta \ell \left( h(x_i, \beta^{(k)}), y_i \right)$$

# Stochastic First-Order Optimization



**Convergence Rate**<sup>4</sup>: worst-case # iterations, in which  $\widehat{R}_D(h_\beta, \ell) \leq \widehat{R}_D(h_{\beta^*}, \ell) + \epsilon$

<sup>4</sup> Bottou, Curtis, and Nocedal, “Optimization Methods for Large-Scale Machine Learning”, 2018.



# Stochastic First-Order Optimization



**Convergence Rate**<sup>4</sup>: worst-case # iterations, in which  $\widehat{R}_D(h_\beta, \ell) \leq \widehat{R}_D(h_{\beta^*}, \ell) + \epsilon$

► **GD**:  $\propto m \cdot \log(\frac{1}{\epsilon})$

► **SGD**:  $\propto \frac{1}{\epsilon}$  (independent of  $m$ )

<sup>4</sup> Bottou, Curtis, and Nocedal, “Optimization Methods for Large-Scale Machine Learning”, 2018.

# Stochastic First-Order Optimization



**Convergence Rate<sup>4</sup>:** worst-case # iterations, in which  $\widehat{R}_D(h_\beta, \ell) \leq \widehat{R}_D(h_{\beta^*}, \ell) + \epsilon$

- ▶ **GD:**  $\propto m \cdot \log(\frac{1}{\epsilon})$
- ▶ **SGD:**  $\propto \frac{1}{\epsilon}$  (independent of  $m$ )
- ▶ For SGD, the same rate applies to  $R(h_\beta, \ell)$  (independent of  $D$  if  $m \gg k$ ) 🚀

Hence, SGD has an amazing performance for large data sets.

<sup>4</sup> Bottou, Curtis, and Nocedal, “Optimization Methods for Large-Scale Machine Learning”, 2018.

# Stochastic First-Order Optimization



**Noise Reduction:** use mini-batches instead of single examples,

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \alpha^{(k)} \frac{1}{b} \sum_{i=1}^b \nabla_{\beta} \ell \left( h(x_{b_i}, \beta^{(k)}), y_{b_i} \right). \quad \text{where } b \ll m.$$

- ▶ smaller variance of update steps
- ▶ stepsize  $\{\alpha^{(k)}\}$  is easier to tune
- ▶ most common approach for deep nets

# Stochastic First-Order Optimization



**Noise Reduction:** use mini-batches instead of single examples,

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \alpha^{(k)} \frac{1}{b} \sum_{i=1}^b \nabla_{\beta} \ell(h(x_{b_i}, \beta^{(k)}), y_{b_i}). \quad \text{where } b \ll m.$$

- ▶ smaller variance of update steps
- ▶ stepsize  $\{\alpha^{(k)}\}$  is easier to tune
- ▶ most common approach for deep nets

## Learning Rate Scheduling:

- ▶ even with mini-batches, noise can eventually prevent the reduction of  $\widehat{R}_D(h_{\beta}, \ell)$
- ▶ hence, decrease step sizes  $\{\alpha^{(k)}\}$  over time

# Stochastic First-Order Optimization



## Momentum:

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - g(\beta^{(k)}) + \gamma^{(k)} \cdot (\beta^{(k)} - \beta^{(k-1)}) \quad \text{where} \quad \begin{cases} g(\beta^{(k)}) & \text{SGD, GD, or mini-batch gradient} \\ \gamma^{(k)} & \text{a weighting parameter} \end{cases}$$

# Stochastic First-Order Optimization



## Momentum:

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - g(\beta^{(k)}) + \gamma^{(k)} \cdot (\beta^{(k)} - \beta^{(k-1)}) \quad \text{where} \quad \begin{cases} g(\beta^{(k)}) & \text{SGD, GD, or mini-batch gradient} \\ \gamma^{(k)} & \text{a weighting parameter} \end{cases}$$

## Accelerated Gradient a.k.a. Nesterov Momentum:

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - g(\beta^{(k)} + \gamma^{(k)} \cdot (\beta^{(k)} - \beta^{(k-1)})) + \gamma^{(k)} \cdot (\beta^{(k)} - \beta^{(k-1)})$$

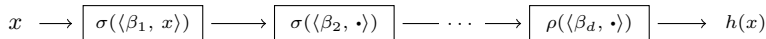
- ▶ momentum is applied before  $g(\cdot)$
- ▶ GD: optimal convergence rate  $\propto \frac{1}{\epsilon^2}$
- ▶ SGD: good practical performance but (theoretical) convergence rate is not improved
- ▶ even better: if combined with adaptive gradients  $\rightarrow$  Adam

# Backpropagation



**Goal:** compute  $\nabla_{\beta} \ell(h(x_i, \beta), y_i)$  where

$$h(x_i, \beta) = \rho\left(\left\langle \beta_d, \phi\left(\left\langle \beta_{d-1}, \dots \phi\left(\left\langle \beta_1, x_i \right\rangle\right)\right\rangle\right)\right\rangle\right)$$



# Backpropagation



**Goal:** compute  $\nabla_{\beta} \ell(h(x_i, \beta), y_i)$  where

$$h(x_i, \beta) = \rho\left(\langle \beta_d, \phi(\langle \beta_{d-1}, \dots \phi(\langle \beta_1, x_i \rangle) \rangle) \rangle\right)$$

$$x \longrightarrow \boxed{\sigma(\langle \beta_1, x \rangle)} \longrightarrow \boxed{\sigma(\langle \beta_2, \cdot \rangle)} \longrightarrow \dots \longrightarrow \boxed{\rho(\langle \beta_d, \cdot \rangle)} \longrightarrow h(x)$$

**Chain rule of calculus:**  $\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(g)}{\partial g} \frac{\partial g(x)}{\partial x}$



# Backpropagation



**Goal:** compute  $\nabla_{\beta} \ell(h(x_i, \beta), y_i)$  where

$$h(x_i, \beta) = \rho\left(\langle \beta_d, \phi(\langle \beta_{d-1}, \dots \phi(\langle \beta_1, x_i \rangle) \rangle) \rangle\right)$$

$$x \longrightarrow \boxed{\sigma(\langle \beta_1, x \rangle)} \longrightarrow \boxed{\sigma(\langle \beta_2, \cdot \rangle)} \longrightarrow \dots \longrightarrow \boxed{\rho(\langle \beta_d, \cdot \rangle)} \longrightarrow h(x)$$

**Chain rule of calculus:**  $\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(g)}{\partial g} \frac{\partial g(x)}{\partial x}$

**Automatic Differentiation:** each function  $f(x)$  also implements its gradient

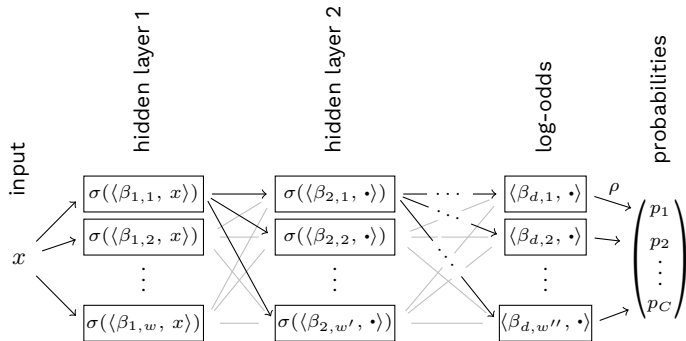
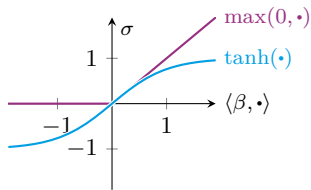
$$\nabla_x f(x) = \left( \frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_n} \right)^\top$$

# Deep Networks



**Deep Nets:** use multiple (logistic regression-like) layers

- ▶ learnable linear combinations  $\langle \beta, \cdot \rangle$
- ▶ non-linear activations  $\sigma$



# Stochastic First-Order Optimization



## Synopsis:

- ▶ **ERM:** we minimize  $\hat{R}_D(h_\beta, \ell) \xrightarrow{m \rightarrow \infty} R(h_\beta, \ell)$
- ▶ **SGD:** gradients randomized through sampling converge quickly for large  $m$

# Stochastic First-Order Optimization



## Synopsis:

- ▶ **ERM:** we minimize  $\widehat{R}_D(h_\beta, \ell) \xrightarrow{m \rightarrow \infty} R(h_\beta, \ell)$
- ▶ **SGD:** gradients randomized through sampling converge quickly for large  $m$
- ▶ **Mini-Batching:** common practice to reduce SGD gradient noise
- ▶ **LR Scheduling:** common practice to balance the noise
- ▶ **Nesterov Momentum:** can improve convergence

# Stochastic First-Order Optimization



## Synopsis:

- ▶ **ERM:** we minimize  $\widehat{R}_D(h_\beta, \ell) \xrightarrow{m \rightarrow \infty} R(h_\beta, \ell)$
- ▶ **SGD:** gradients randomized through sampling converge quickly for large  $m$
- ▶ **Mini-Batching:** common practice to reduce SGD gradient noise
- ▶ **LR Scheduling:** common practice to balance the noise
- ▶ **Nesterov Momentum:** can improve convergence

## Practical Recommendations:

- ▶ **Carefully Design Loss Functions** to reflect your goals

# Stochastic First-Order Optimization



## Synopsis:

- ▶ **ERM:** we minimize  $\widehat{R}_D(h_\beta, \ell) \xrightarrow{m \rightarrow \infty} R(h_\beta, \ell)$
- ▶ **SGD:** gradients randomized through sampling converge quickly for large  $m$
- ▶ **Mini-Batching:** common practice to reduce SGD gradient noise
- ▶ **LR Scheduling:** common practice to balance the noise
- ▶ **Nesterov Momentum:** can improve convergence

## Practical Recommendations:

- ▶ **Carefully Design Loss Functions** to reflect your goals
- ▶ **Use Popular First-Order Methods** like Adam or SGD with Nesterov Momentum



# Data and Assumptions

## A Premature Conclusion



Machine learning = data ○ model ○ fit

### What we have learned:

- ▶ Deep Nets are universal function approximators
- ▶ Customized loss functions let them learn what we need
- ▶ We know effective ways of optimizing them



# A Premature Conclusion



**Machine learning** = data ○ **model** ○ fit

## What we have learned:

- ▶ Deep Nets are universal function approximators
- ▶ Customized loss functions let them learn what we need
- ▶ We know effective ways of optimizing them

**What could possibly go wrong?** 💣

# Learning Assumptions



Recall that we approximate

$$R(h_\beta, \ell) = \mathbb{E}_{(x,y) \sim \mathbb{P}}(\ell(h_\beta(x), y))$$

through

$$\hat{R}_D(h_\beta, \ell) = \frac{1}{m} \sum_{i=1}^m \ell(h_\beta(x_i), y_i)$$

# Learning Assumptions



Recall that we approximate

$$R(h_\beta, \ell) = \mathbb{E}_{(x,y) \sim \mathbb{P}}(\ell(h_\beta(x), y))$$

through

$$\hat{R}_D(h_\beta, \ell) = \frac{1}{m} \sum_{i=1}^m \ell(h_\beta(x_i), y_i)$$

**Independent and Identical Distribution (IID) Assumption:**

$$(x, y) \sim \mathbb{P} \quad \forall (x, y) \in D \cup D_{\text{test}}$$

# Learning Assumptions



Recall that we approximate

$$R(h_\beta, \ell) = \mathbb{E}_{(x,y) \sim \mathbb{P}}(\ell(h_\beta(x), y))$$

through

**Data Set Shift** breaks the IID assumption 💣

- ▶  $D \sim \mathbb{P}_S$  (e.g., a *simulation*)
- ▶  $D_{\text{test}} \sim \mathbb{P}_T$  (e.g., a *real detector*)
- ▶  $\mathbb{P}_S \neq \mathbb{P}_T$

$$(x, y) \sim \mathbb{P} \quad \forall (x, y) \in D \cup D_{\text{test}}$$

## Types of Data Set Shift<sup>5</sup>



Recognize that  $\mathbb{P}(X, Y) = \mathbb{P}(X | Y) \cdot \mathbb{P}(Y)$

<sup>5</sup> Kull and Flach, “Patterns of dataset shift”, 2014

## Types of Data Set Shift<sup>5</sup>



Recognize that  $\mathbb{P}(X, Y) = \mathbb{P}(X | Y) \cdot \mathbb{P}(Y)$

### Label Shift:

$$\mathbb{P}_{\mathcal{S}}(X | Y) = \mathbb{P}_{\mathcal{T}}(X | Y)$$

$$\mathbb{P}_{\mathcal{S}}(Y) \neq \mathbb{P}_{\mathcal{T}}(Y)$$

<sup>5</sup> Kull and Flach, “Patterns of dataset shift”, 2014

# Types of Data Set Shift<sup>5</sup>



Recognize that  $\mathbb{P}(X, Y) = \mathbb{P}(X | Y) \cdot \mathbb{P}(Y)$

## Label Shift:

$$\mathbb{P}_{\mathcal{S}}(X | Y) = \mathbb{P}_{\mathcal{T}}(X | Y)$$

$$\mathbb{P}_{\mathcal{S}}(Y) \neq \mathbb{P}_{\mathcal{T}}(Y)$$

## Concept Shift:

$$\mathbb{P}_{\mathcal{S}}(Y) = \mathbb{P}_{\mathcal{T}}(Y)$$

$$\mathbb{P}_{\mathcal{S}}(X | Y) \neq \mathbb{P}_{\mathcal{T}}(X | Y)$$

<sup>5</sup> Kull and Flach, “Patterns of dataset shift”, 2014

## Types of Data Set Shift<sup>5</sup>



Recognize that  $\mathbb{P}(X, Y) = \mathbb{P}(X | Y) \cdot \mathbb{P}(Y)$

$$= \mathbb{P}(X) \cdot \mathbb{P}(Y | X)$$

### Label Shift:

$$\mathbb{P}_{\mathcal{S}}(X | Y) = \mathbb{P}_{\mathcal{T}}(X | Y)$$

$$\mathbb{P}_{\mathcal{S}}(Y) \neq \mathbb{P}_{\mathcal{T}}(Y)$$

### Concept Shift:

$$\mathbb{P}_{\mathcal{S}}(Y) = \mathbb{P}_{\mathcal{T}}(Y)$$

$$\mathbb{P}_{\mathcal{S}}(X | Y) \neq \mathbb{P}_{\mathcal{T}}(X | Y)$$

<sup>5</sup> Kull and Flach, “Patterns of dataset shift”, 2014



# Types of Data Set Shift<sup>5</sup>



Recognize that  $\mathbb{P}(X, Y) = \mathbb{P}(X | Y) \cdot \mathbb{P}(Y)$

$$= \mathbb{P}(X) \cdot \mathbb{P}(Y | X)$$

## Label Shift:

$$\mathbb{P}_{\mathcal{S}}(X | Y) = \mathbb{P}_{\mathcal{T}}(X | Y)$$

$$\mathbb{P}_{\mathcal{S}}(Y) \neq \mathbb{P}_{\mathcal{T}}(Y)$$

## Concept Shift:

$$\mathbb{P}_{\mathcal{S}}(Y) = \mathbb{P}_{\mathcal{T}}(Y)$$

$$\mathbb{P}_{\mathcal{S}}(X | Y) \neq \mathbb{P}_{\mathcal{T}}(X | Y)$$

## (Also) Concept Shift:

$$\mathbb{P}_{\mathcal{S}}(X) = \mathbb{P}_{\mathcal{T}}(X)$$

$$\mathbb{P}_{\mathcal{S}}(Y | X) \neq \mathbb{P}_{\mathcal{T}}(Y | X)$$

<sup>5</sup> Kull and Flach, “Patterns of dataset shift”, 2014

# Types of Data Set Shift<sup>5</sup>



Recognize that  $\mathbb{P}(X, Y) = \mathbb{P}(X | Y) \cdot \mathbb{P}(Y)$

$$= \mathbb{P}(X) \cdot \mathbb{P}(Y | X)$$

## Label Shift:

$$\mathbb{P}_{\mathcal{S}}(X | Y) = \mathbb{P}_{\mathcal{T}}(X | Y)$$

$$\mathbb{P}_{\mathcal{S}}(Y) \neq \mathbb{P}_{\mathcal{T}}(Y)$$

## Concept Shift:

$$\mathbb{P}_{\mathcal{S}}(Y) = \mathbb{P}_{\mathcal{T}}(Y)$$

$$\mathbb{P}_{\mathcal{S}}(X | Y) \neq \mathbb{P}_{\mathcal{T}}(X | Y)$$

## (Also) Concept Shift:


$$\mathbb{P}_{\mathcal{S}}(X) = \mathbb{P}_{\mathcal{T}}(X)$$

$$\mathbb{P}_{\mathcal{S}}(Y | X) \neq \mathbb{P}_{\mathcal{T}}(Y | X)$$

## Covariate Shift:

$$\mathbb{P}_{\mathcal{S}}(Y | X) = \mathbb{P}_{\mathcal{T}}(Y | X)$$

$$\mathbb{P}_{\mathcal{S}}(X) \neq \mathbb{P}_{\mathcal{T}}(X)$$

**Correction Methods** are available for each type, but require extra information  
(additional data, more assumptions, ...) 

<sup>5</sup> Kull and Flach, “Patterns of dataset shift”, 2014

# Domain-Adversarial Unsupervised Domain Adaptation<sup>6</sup>



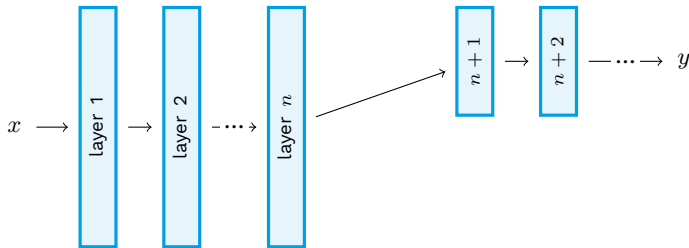
- ▶ **Assume Concept Shift**  $\mathbb{P}_{\mathcal{S}}(X | Y) \neq \mathbb{P}_{\mathcal{T}}(X | Y)$  and  $\mathbb{P}_{\mathcal{S}}(Y) = \mathbb{P}_{\mathcal{T}}(Y)$
- ▶ **Employ Unlabeled Data**  $D_{\mathcal{T}} = \{x \sim \mathbb{P}_{\mathcal{T}}(X)\}$

<sup>6</sup> Ganin et al., “Domain-Adversarial Training of Neural Networks”, 2016

# Domain-Adversarial Unsupervised Domain Adaptation<sup>6</sup>



- ▶ **Assume Concept Shift**  $\mathbb{P}_{\mathcal{S}}(X | Y) \neq \mathbb{P}_{\mathcal{T}}(X | Y)$  and  $\mathbb{P}_{\mathcal{S}}(Y) = \mathbb{P}_{\mathcal{T}}(Y)$
- ▶ **Employ Unlabeled Data**  $\mathcal{D}_{\mathcal{T}} = \{x \sim \mathbb{P}_{\mathcal{T}}(X)\}$

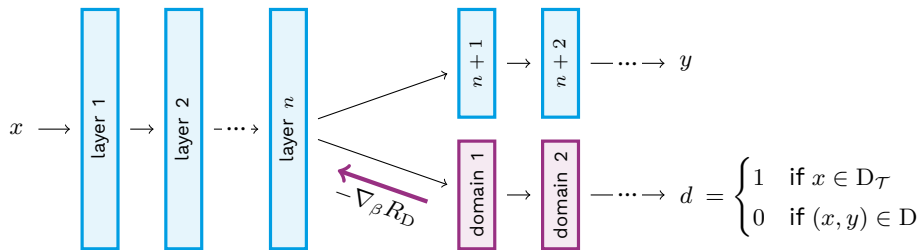


<sup>6</sup> Ganin et al., “Domain-Adversarial Training of Neural Networks”, 2016

# Domain-Adversarial Unsupervised Domain Adaptation<sup>6</sup>



- **Assume Concept Shift**  $\mathbb{P}_{\mathcal{S}}(X | Y) \neq \mathbb{P}_{\mathcal{T}}(X | Y)$  and  $\mathbb{P}_{\mathcal{S}}(Y) = \mathbb{P}_{\mathcal{T}}(Y)$
- **Employ Unlabeled Data**  $D_{\mathcal{T}} = \{x \sim \mathbb{P}_{\mathcal{T}}(X)\}$



<sup>6</sup> Ganin et al., “Domain-Adversarial Training of Neural Networks”, 2016

# Class-Conditional Label Noise<sup>7</sup>



## Label Noise:

- ▶ **Training Labels**  $\hat{y}$  are randomly flipped versions of the ground-truth  $y$
- ▶ **Assumptions** about the flipping process  $y \rightarrow \hat{y}$  are required

<sup>7</sup> Menon et al., “Learning from Corrupted Binary Labels via Class-Probability Estimation”, 2015

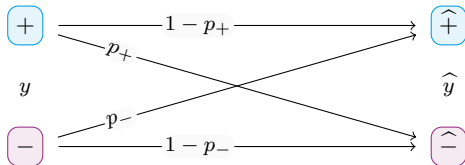
## Class-Conditional Label Noise<sup>7</sup>



### Label Noise:

- ▶ **Training Labels**  $\hat{y}$  are randomly flipped versions of the ground-truth  $y$
- ▶ **Assumptions** about the flipping process  $y \rightarrow \hat{y}$  are required

**Class-Conditional Noise:**  $\mathbb{P}(Y = +1 \mid X = x) = a \cdot \mathbb{P}(\hat{Y} = +1 \mid X = x) + b$



<sup>7</sup> Menon et al., “Learning from Corrupted Binary Labels via Class-Probability Estimation”, 2015

## Deep Sets<sup>8</sup>



- ▶ **Each instance is a set**  $\{x_i \in \mathcal{X} : 1 \leq i \leq m\}$  of variable size  $m$
- ▶  $\mathcal{Y}$  are properties of such sets

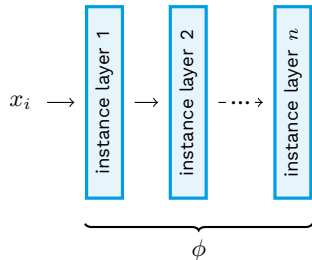
<sup>8</sup> Zaheer et al., “Deep sets”, 2017



## Deep Sets<sup>8</sup>



- ▶ **Each instance is a set**  $\{x_i \in \mathcal{X} : 1 \leq i \leq m\}$  of variable size  $m$
- ▶  $\mathcal{Y}$  are properties of such sets

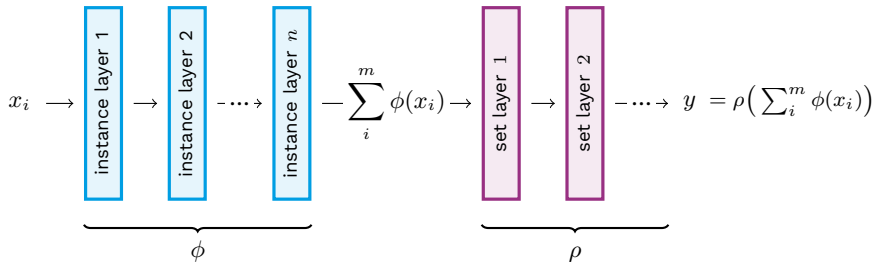


<sup>8</sup> Zaheer et al., “Deep sets”, 2017

## Deep Sets<sup>8</sup>



- ▶ **Each instance is a set**  $\{x_i \in \mathcal{X} : 1 \leq i \leq m\}$  of variable size  $m$
- ▶  $\mathcal{Y}$  are properties of such sets

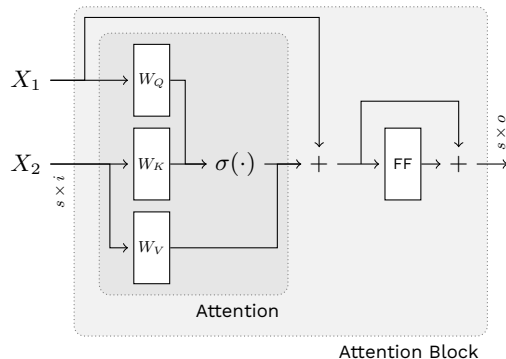
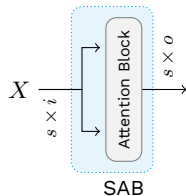


<sup>8</sup> Zaheer et al., “Deep sets”, 2017

# Set Transformers<sup>9</sup>



**Self-Attention Block:**

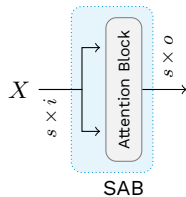


<sup>9</sup> Lee et al., “Set Transformer: A Framework for Attention-based Permutation-Invariant Neural Networks”, 2019

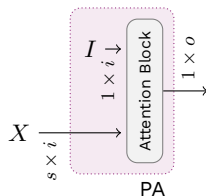
# Set Transformers<sup>9</sup>



**Self-Attention Block:**

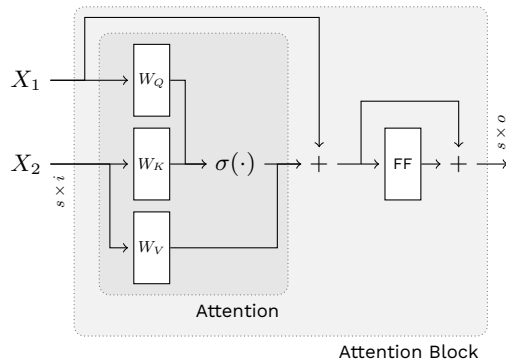


**Pooling by Attention:**



**Set Encoder:**

$\text{SAB} \circ \dots \circ \text{SAB} \circ \text{PA}$



<sup>9</sup> Lee et al., “Set Transformer: A Framework for Attention-based Permutation-Invariant Neural Networks”, 2019



## Concluding Remarks

# Should I Use Neural Networks?



**Architecture Search** vs feature engineering

**Scale** great for big data (but not for small data)

**GPUs required** as well as computation time for fitting

# Implementing Neural Networks



## JAX, PyTorch, Tensorflow, or Keras?

- ▶ Keras, Tensorflow: established solutions
- ▶ PyTorch, JAX: maximum flexibility

# Implementing Neural Networks



## JAX, PyTorch, Tensorflow, or Keras?

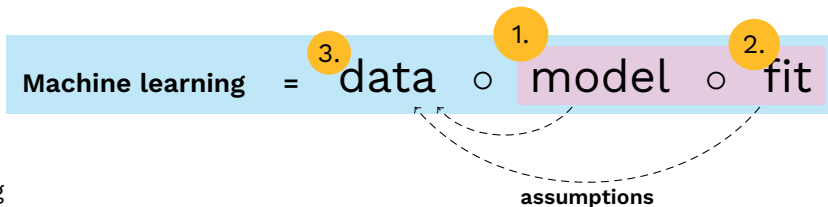
- ▶ Keras, Tensorflow: established solutions
- ▶ PyTorch, JAX: maximum flexibility

### JAX:

- ▶ JIT compilation speedups
- ▶ API identical to Numpy/Scipy
- ▶ Clean functional programming style (clarity, separation of concerns)
- ▶ Evolving eco-system and fewer solutions



# Agenda



1. Modeling
2. Fitting
3. Data and Assumptions
4. Concluding Remarks

+ Hands-On Exercises (Tue ~ 45 min, Thu ~ 90 min)



`https:  
//git.e5.physik.tu-dortmund.de/qfuehring/ml_intro_handson`