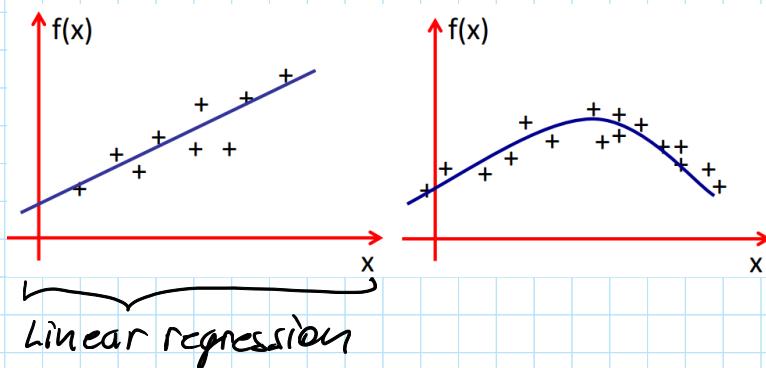


## Linear Regression

- Input:  $x$ : Age, Sex, BMI
- Output:  $y$ : Measure of disease progression
- Goal: Learn real valued mapping



- Linear regression:  $\cdot y \approx f(x)$

$$f(x) = \omega^T x + b \quad (\text{e.g. } ax + b)$$

$$\hat{\omega}^T \hat{x}$$

$$\cdot \hat{x} = [x_1, x_2, \dots, x_d, 1]$$

$$\cdot \hat{\omega} = [\omega_1, \omega_2, \dots, \omega_d, b]$$

• How to find optimal weight vector:

$$\hookrightarrow \hat{\omega}^* = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \omega^T x_i)^2$$

$\hookrightarrow$  Least squares error  
 $\hookrightarrow$  Minimize error

- Residual:  $\cdot r_i = y_i - f(x_i)$

$$\hookrightarrow \text{"Restwert"} = y_i - \omega^T x_i$$

- Squared Residual:  $R(\omega) = \sum_{i=1}^n (y_i - \omega^T x_i)^2$

- Closed form solution:  $\cdot \hat{\omega}^* = (X^T X)^{-1} \cdot X^T y$

$\hookrightarrow x$  = Data point inputs  
 $\hookrightarrow y$  = Outputs corresponding

$\hookrightarrow$  to Inputs  
 $\hookrightarrow$  to Inputs

$\cdot$  computational complexity:  $O(d^3)$

$\cdot$  Input data is just reference data  
 $\hookrightarrow$  No need for exact solution

- Optimization:
  - Minimize  $\hat{R}(w) = \sum_i (y_i - w^T x_i)^2$
  - with gradient descend
  - $\hat{R}(w)$  is convex: local minimum = global minimum
  - 1. Start at an arbitrary point  $w_0 \in \mathbb{R}^d$
  - 2.  $w_{t+1} = w_t - \eta \nabla \hat{R}(w_t)$
  - $\eta$ : - Learning rate / step size
    - If we choose a poor step size  
 $\hookrightarrow$  oscillation 

### Adaptive step size:

- Line search: -  $\eta_t = \operatorname{argmin} \hat{R}(w_t - \eta \cdot \nabla \hat{R}(w_t))$ 
  - choose stepsize that gives biggest improvement towards 0
- Bold driver:
  - If function decreases, increase step-size  
 $\hookrightarrow \hat{R}(w_{t+1}) < \hat{R}(w_t) \Rightarrow \eta_{t+1} = \eta_t \cdot 1,1$
  - If function increases, decrease step-size  
 $\hookrightarrow \hat{R}(w_{t+1}) > \hat{R}(w_t) \Rightarrow \eta_{t+1} = \eta_t \cdot 0,5$

### Loss functions

- Least squared loss:  $l(r) = r^2$ 
  - $\hookrightarrow$  If I under/overpredict, it's the same, because square gets rid of sign
- Absolute loss:  $l(r) = |r|$ 
  - $\hookrightarrow$  Overvalue small errors / undervalue big errors
  - $\hookrightarrow$  More robust to ausgeschlagen

## Over-/Underfitting

- Overfitting: Too high degree polynomial
- Underfitting: Too low degree polynomial
- The higher the degree, the better we fit the data

## Prediction Error

- Wenn wir zufällig ein  $x$  einfügen, was ist dann der expected error
- We assume, our testdata is drawn from the distribution  $P$ :  $(x_i, y_i) \sim P(x, y)$
- True risk:
  - Expected error
  - $R(w) = \int P(x, y) \cdot (y - w^T x)^2 \cdot dx \cdot dy$
  - $= E_{x,y} [(y - w^T x)^2]$
  - cannot be computed since we don't know the real distribution of  $P$

## Estimating true risk

- Empirical risk  $\hat{R}_D(w)$

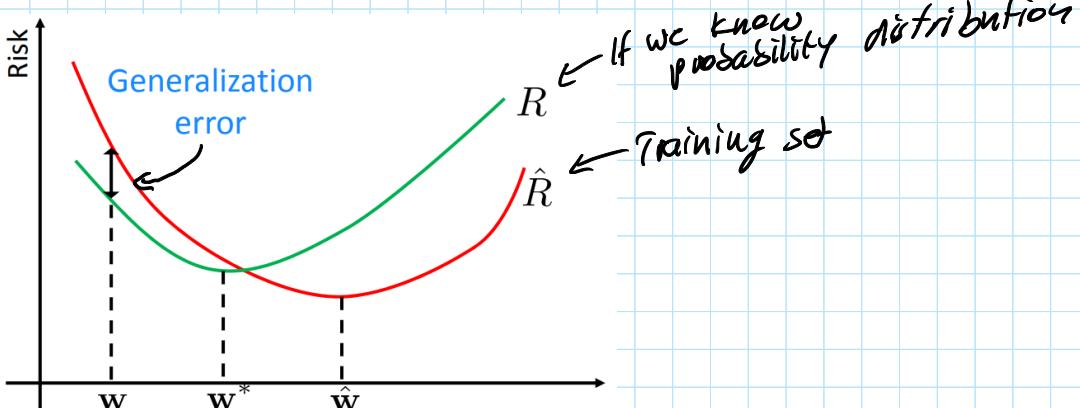
↳ If we have enough data points

↳ Law of large numbers

$$\hat{R}_D(w) = \frac{1}{|D|} \cdot \sum_{\substack{(x_i, y_i) \\ \in D}} (y_i - w^T x_i)^2$$

Total error

↳ Was ist der expected error, wenn wir  $y$  ausgeben



To be successful:  $|R(w) - \hat{R}_D(w)| \rightarrow 0$  as  $D \rightarrow \infty$

$$\cdot E[R_{\text{trainData}}(\hat{w})] \quad \underbrace{\quad}_{\text{only training data}}$$

$$\ll E[R(\hat{w})] \quad \underbrace{\quad}_{\text{whole distribution}}$$

Da wir viel mehr haben

$\rightarrow$  Underestimating prediction error

### Avoid underestimating prediction error

- Split data into training and test set

1. Optimize  $\hat{w}$  using training data

$$\hat{w} = \underset{w}{\operatorname{argmin}} \hat{R}_{\text{trainData}}(w)$$

2. Evaluate prediction error on test data

$$\hat{R}_{\text{testData}}(\hat{w}) = \frac{1}{|D_{\text{testData}}|} \cdot \sum_{\substack{(x,y) \in \\ \text{Test data}}} (y - \hat{w}^T x)^2$$

### Choosing polynomial degree

1. Do for each degree  $m$  the following for  $i=1, \dots, k$

2. Split data into training / test-data:  $D = D_{\text{Train}}^{(i)} \cup D_{\text{test}}^{(i)}$

3. Train model:  $\hat{w}_{i,m} = \underset{w}{\operatorname{argmin}} \hat{R}_{\text{Train}}^{(i)}(w)$

4. Estimate error:  $\hat{R}_m^{(i)} = \hat{R}_{\text{test}}^{(i)}(\hat{w}_{i,m})$

5. Select polynomial  $\hat{m} = \underset{m}{\operatorname{argmin}} \frac{1}{k} \sum_{i=1}^k \hat{R}_m^{(i)}$

$\underbrace{\quad}_{\text{Durchschnittlicher Fehler für Polynom mit grad } m}$

### Cross validation

- Monte-Carlo cross validation: Randomly

-  $k$ -fold cross validation: 1. Partition data into  $k$ -folds

2. Train on  $(k-1)$  folds

3. Test on 1 fold

- How to choose # folds?

- Too small:
  - Overfitting to test data
  - Underfitting to training data

Regularization (- Encourage small weights with penalty)

- If we only optimize loss functions, we'll get large weights

Ridge Regression (= L2-regularized Regression)

$$-\hat{w} = \min \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \underbrace{\|w\|_2^2}_{\sum w_i^2}$$

$$-\text{closed solution: } \hat{w} = (X^T X + \lambda I)^{-1} \cdot X^T y$$

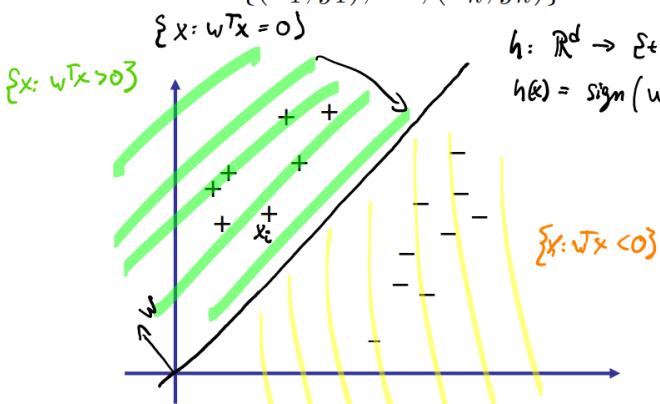
- How to choose  $\lambda$ : cross validation

- The greater  $\lambda$ , the lower the weights

Classification

-  $y$  is discrete (e.g. spam/not spam)

• Data set  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$



-  $w^T x$  either  $> 0 / < 0$

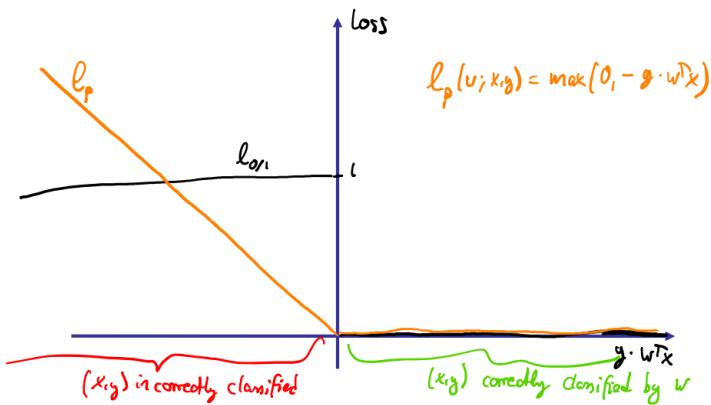
-  $h(x) = \text{sign}(w^T x)$

- We have to find a linear separator

- 0/1 loss is not convex  $\Rightarrow$  Surrogate loss

## Surrogate Loss

- Solve  $\hat{w} = \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l_p(w, x_i, y_i)$
- Perceptron loss:  $l_p(w, x_i, y_i) = \max(0, -y_i \cdot w^T x_i)$
- $-y_i \cdot w^T x_i > 1$  if  $y_i = +1$  and  $w^T x_i = -1$   
 $y_i = -1$  and  $w^T x_i = +1$  } Error

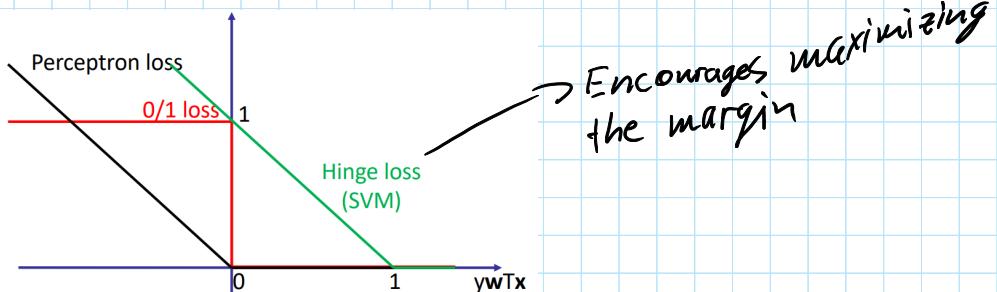


## Stochastic gradient descent

- Computing the gradient requires summing over all data
- For large data sets, this is inefficient
- Use a part of the set or even only one point
- Mini batches: Averaging over the gradients of randomly selected points

## Support Vector Machines (=SVMs)

- Want to maximize the margin/boundary between separating line and first data points
- Hinge Loss:  $l_H(w, x_i, y) = \max(0, 1 - y \cdot w^T x)$



$$\hat{\omega} = \underset{\omega}{\operatorname{arg\,min}} \frac{1}{n} \cdot \sum_{i=1}^n \max(0, 1 - y_i \cdot \omega^\top x_i) + \lambda \underbrace{\|\omega\|_2^2}_{\text{Regularization}}$$

## Feature Selection

### - Greedy feature selection

- Greedily add/remove features to maximize

• Set of all features  $X = \{x_1, \dots, x_d\}$

1. Only train with subset  $S$

2.  $\hat{L}(S) :=$  cross validation error using only features in  $S$

### - Greedy forward selection

• for  $i = 1, \dots, d$

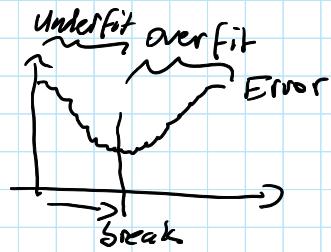
• find best feature to add:  $s_i = \underset{j \in V \setminus S}{\operatorname{arg\,min}} \hat{L}(S \cup \{x_j\})$

• Compute error with  $S \cup s_i$ :  $E_i$  underfit overfit

• If  $E_i > E_{i-1} \Rightarrow$  break

• Else  $S \leftarrow S \cup \{x_i\}$

• Faster than backward



### - Greedy backward selection

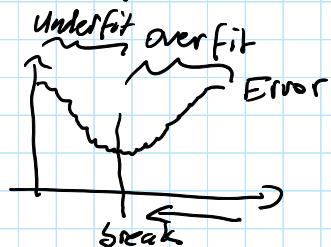
• for  $i = d-1, \dots, 1$

• find best feature to remove:  $s_i = \underset{j \in S}{\operatorname{arg\,min}} \hat{L}(S \setminus \{x_j\})$

• Compute error without  $s_i$ :  $E_i$  underfit overfit

• If  $E_i > E_{i+1}$  break

• Else  $S \leftarrow S \setminus \{x_i\}$



## Sparse Regression: Lasso (=L1-regularized regression)

- We only want the most important features, therefore constraining the sparsity of  $w$  induces how many features are used
- Replace  $\|w\|_2^2$  by  $\|w\|_1$
- min  $\sum_w (y_i - w^T x_i) + \lambda \cdot \|w\|_1$
- This encourages coefficients to be exactly 0 which induces automatic feature selection
- Can be used in SVM, too!

## Kernels

- Used for non-linear classification
- The weights are just:  $w = \sum y_i \cdot \alpha_i \cdot x_i$
- We only need the inner / dot product  
 $\hookrightarrow x^T x \rightarrow \phi(x)^T \phi(x') =: k(x, x')$
- Polynomial kernel:  $k(x, x') = (1 + x^T x')^m$
- Kernel Trick:
  1. Express problem such it only depends on inner product  $x^T \cdot x$
  2. Replace inner product by kernel
- Reformulating perceptron:  
 $\hat{w} = \arg \min \frac{1}{n} \sum \max(0, -y_i \cdot w^T x_i)$   $\text{---}^*$
- Kernel requirements:
  - Symmetric  $k(x', x) = k(x, x')$
  - Positive semi-definite
$$\begin{aligned} & -y_i (\sum y_j \cdot \alpha_j \cdot x_j)^T x_i \\ & -y_i \sum y_j \cdot \alpha_j \cdot \underbrace{(x_j^T \cdot x_i)}_{\text{inner product}} \end{aligned}$$
- Kernel engineering:
  - $c \cdot k_1(x, x') + [k_2(x, x') \cdot k_3(x, x')] + f(k_4(x, x'))$
  - Composition of kernels are valid kernels

## - Kernels as similarity function

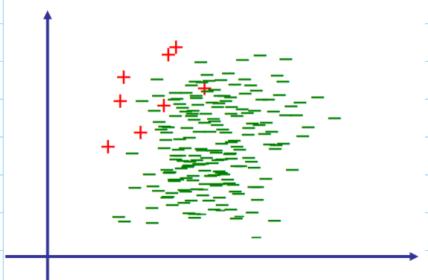
- Given two D dimensional vectors  $x_i$  and  $x_j$
- The Gaussian kernel is defined as  $k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{\sigma^2}\right)$
- $\|x_i - x_j\|^2$  is the euclidean distance between  $x_i$  and  $x_j$
- so the Gaussian kernel is a similarity measure
- If it evaluates to 1  $\Leftrightarrow x_i = x_j$
- And the more it approaches 0, the further apart they are

## Parametric vs. non-Parametric learning

- Parametric models:
  - finite set of data according to dimension
  - Even if we add new  $(x, y)$  pairs
    - e.g. linear regression
- Nonparametric model:
  - Model grows in complexity with the size of data
  - Kernahized perceptron
- Kernels convert parametric to nonparametric
  - for each new  $(x, y)$  pair, we need new  $\alpha_j$  variable

## Class Imbalance

- Sources of imbalanced data
  - Fraud detection
  - Spam filtering
- Subsampling: Remove data from majority class randomly
- Upsampling: Repeat data from minority class with small perturbations



## - Cost sensitive classification

- Modify Perceptron / SVM to take class balance into consideration
- Add variable  $c$  to loss function which has different values for  $+$ / $-$  controlling the tradeoff

- Perceptron:  $\ell_{CS-P}(\mathbf{w}; \mathbf{x}, y) = c_y \max(0, -y\mathbf{w}^T \mathbf{x})$
- SVM:  $\ell_{CS-H}(\mathbf{w}; \mathbf{x}, y) = c_y \max(0, 1 - y\mathbf{w}^T \mathbf{x})$

with parameters  $c_+, c_- > 0$  controlling tradeoff

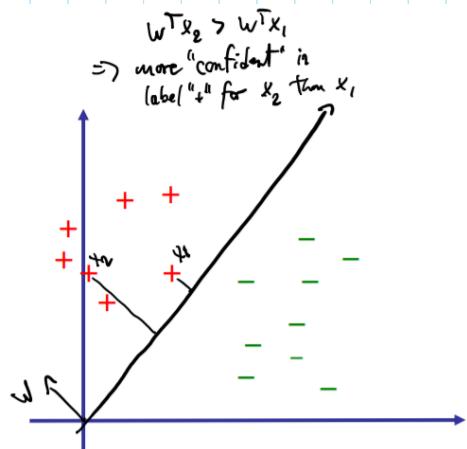
		True label	
		Positive	Negative
Predicted label	Positive	TP	FP
	Negative	FN	TN

$P_+ + P_- = n_+ + n_- = n$        $\Sigma = n_+$        $\Sigma = n_-$

$\Sigma > P_+$   
 $\Sigma = P_-$

## Dealing with multiple classes

- One vs. all: solve  $C$  binary classifiers, one for each class
  - Positive: all points from class  $C$
  - Negative: all points from other classes
  - classify  $x$  by looking for classifier with largest confidence



} Je weiter weg der Punkt von der Linie ist, desto höher ist der Wert von  $w^T x$

- One vs. one:
  - Use a classifier for each pair of classes  $(i, j)$
  - Class with highest number of positive predictions, win
- Multi-class-hinge loss:
  - Have a weight-vector for each of the  $c$  classes
  - $w^{(1)}, \dots, w^{(c)}$
  - Given a data-point, we want to train
 
$$\underbrace{w^{(v)^T} x}_{\substack{\text{confidence} \\ \text{in correct} \\ \text{label}}} \geq \underbrace{w^{(i)^T} x + 1}_{\substack{\text{confidence} \\ \text{in other labels}}}$$

## Artificial Neural Networks (ANNs/Multi layer perceptron)

- Each neuron holds a value:
  - Between 0-1
  - $v_i^{(j)} = \varphi(w_i^{(j)^T} \cdot x)$
- Activation function:
  - $\varphi(x) \rightarrow (0, 1)$
  - e.g. sigmoid function
- Training:
  - We use backward propagation
  - Randomly set all weight values
  - Apply loss function to output
  - When predicting multiple outputs at the same time, define loss as sum of per output losses
  - Optimize the weights to minimize loss
 
$$\hookrightarrow w^* = \arg \min \sum \ell(w; y_i, x_i)$$
  - Use stochastic gradient descent

- Prediction: Using forward propagation

- Predict  $y_i = f_i$  for regression } 1 output
- Predict  $y_i = \text{sign}(f_i)$  for classification
- Use maximum value of all outputs

- Overfitting:  
- Don't run SGD until convergence

- Add regularization to keep weight small

## Convolutional Neural Network (=CNN)

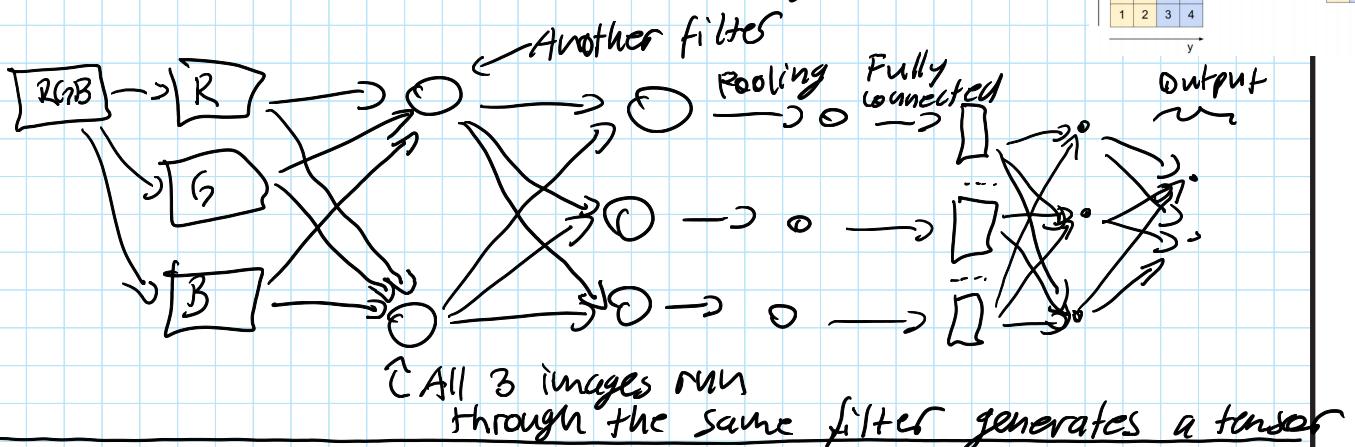
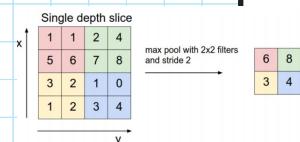
- Input  $\rightarrow$  Convolution  $\rightarrow$  Pooling  $\rightarrow$  Fully connected  $\rightarrow$  Output

- Filter:  
- Padding on borders to let each pixel be at the center

- Stride: move  $(x, x)$  at a time (normally  $(1, 1)$ )  
↳ helps reducing the size of output

- Filters are initialized randomly

- Pooling:  
- Aggregate several units to decrease width  
of the network by averaging them



## Learning with momentum

- Can help to escape a local minimum

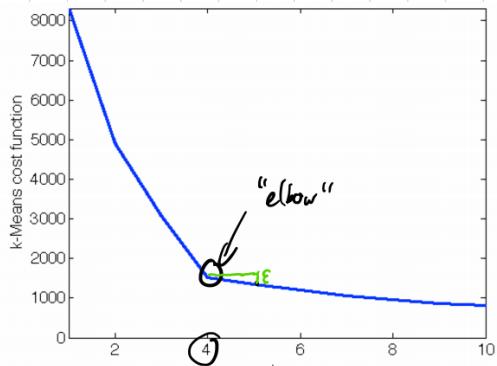
- Move not only in the direction of the gradient,  
but also in the direction of the last weight update

$$a = m \cdot a + \eta \cdot \nabla \ell(w, y, x)$$

↑ momentum < 1

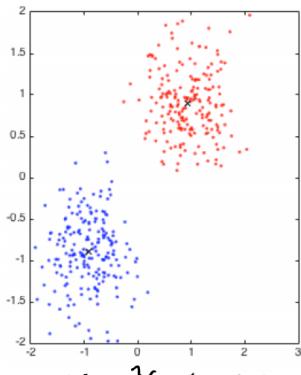
## Clustering

- Given data points, group into clusters
  - K-means:
    - Represent each cluster by single center point
    - Assign points to closest center-point
1. Initialize cluster centers  $\mu^{(0)} = [\mu_1^{(0)}, \dots]$  randomly
  2. Assign each point to closest center
  3. Update center as mean of all points assigned to that center
- How to choose # of classes = centers:

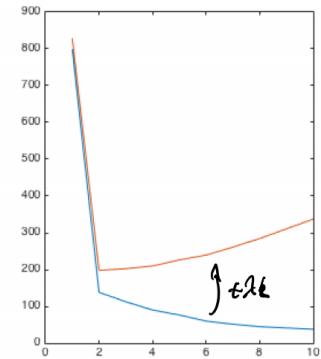


- Can't use cross validation because loss keeps decreasing  
↳ Look for elbow

- you can add regularization to add penalty for each additional class



$$\min_{k, \mu_{1:k}} \hat{R}(\mu_{1:k}) + \lambda \cdot k$$



Equivalent to "elbow method" for  $\epsilon = 2$

- K-means++:
    - Start with one center randomly
    - Add centers  $2, \dots, k$  randomly, proportionally to squared distance to closest selected center
- 

### Dimension Reduction

- We want to reduce dimension from  $d$  to  $k$
- We have to translate features  $x_{(d, i)}$  to  $z_{(k, i)}$
- $w \cdot z_i = x_i$  (we have to find  $w$  with dimensions  $(d \times k)$ )
- $z_i = w^T x_i$
- To get unique results, we have to center data to covariance-matrix
 
$$\hookrightarrow \Sigma = \frac{1}{n} \sum_{i=1}^d x_i \cdot x_i^T = \frac{1}{n} \cdot X^T \cdot X$$
- PCA:
  - $(w, z_1, \dots, z_n) = \underset{\text{argmin}}{\text{argmin}} \sum_{i=1}^n \|w z_i - x_i\|_2^2$
  - $w$  is either given by Eigendecomposition or SVD
  - $w = (v_1 | \dots | v_k)$ 

$$\Sigma = \underbrace{\sum_{i=1}^d \lambda_i v_i \cdot v_i^T}_{\Sigma = U \Sigma V^T}$$

- Kernel PCA:
  - For non-linear dimension reduction

- First increase dimension to get linear mapped data and then reduce dimension to  $k$

- $K = K^T \cdot K$

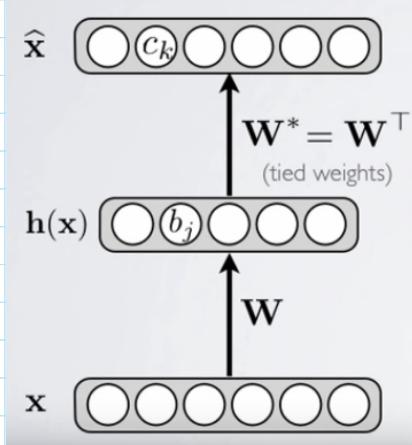
- $\alpha^* = \arg \max \alpha^T \cdot K^T \cdot K \cdot \alpha$

- $\alpha^{(1)} = \frac{1}{\sqrt{\lambda_1}} \cdot v_1$

- A new point is projected:

$$\hookrightarrow z_i = w^T x = \underbrace{\left( \sum \alpha_j \cdot x_j \right)^T}_{w} \cdot x = \sum \alpha_j \cdot (x_j^T x) = \sum \alpha_j \cdot k(x_j, x)$$

- Autoencoder:



-  $x \mapsto \hat{x}$

- we use  $W$  as

heights for  
dimension reduction

### Probability Model

- What is the best prediction if we would know the probability and could use any function  $h(x)$
- Least squares:
  - $\min_h E_{x,y} [(y - h(x))^2]$
  - $= E_x \left[ \min_{h(x)} E_y [(y - h(x))^2 | X=x] \right]$ 
    - ↑ For each  $x$  we could use another function
  - $\hat{y} = E[Y | X=x]$ 
    - ↳ Conditional mean is the best predictor for  $y$ , given  $x$
    - ↳ Baye's optimal predictor
  - But we never have true distribution, therefore estimate conditional distribution
    - ↳ use MLE to search for parameters for distribution with using the train-data
  - When solving linear regression, we do maximum likelihood on the data

## Logistic Regression (= Discriminative)

- predicts whether something is true/false instead of a continuous output
- Returns a probability between (0,1)
- We use Maximum Likelihood for positioning the S-curve to the data
- Use S-curve to get probability

## Bayesian decision theory

- Recommendation for picking the action that minimizes the expected cost
- $a^* = \underset{a \in A}{\operatorname{argmin}} E_y[C(y, a)](x)$
- The action that minimizes the cost is the most likely function
- Assymmetric cost: if we have more costs for e.g. FP than FN
- Doubtful: pick most likely class only if confident enough  
 $\hookrightarrow a^* = \begin{cases} y & \text{if } P(y|x) \geq 1 - c \\ D & \text{otherwise} \quad (= \text{Don't know}) \end{cases}$

→ labeling features,  
e.g. Σ-1.13

## Uncertainty sampling

- Labels (test/train-data) is expensive
- Always pick the example that we are most uncertain about

1. Estimate  $\hat{P}(y|x)$  for all labels

2. Pick most uncertain label:  $\arg\min_i |\theta_j - \hat{P}(y_j|x_i)|$

---

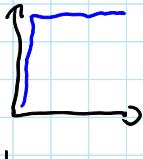
## Generalization error

- Difference between the expected and empirical loss
  - Expected loss: Expected error over all possible values of  $(x,y)$  when we know the true distribution
  - Empirical loss: Empirical error on sample data  
 $\hookrightarrow \frac{1}{n} \sum R(w, x, y)$
- 

## ROC and Precision-Recall-curves

### - ROC:

- X-axis: false positives
- Y-axis: true positives
- Used for balanced classes



### - PR:

- X-axis: recall
- Y-axis: precision
- Used when there are class imbalances



## Metrics for classification

- True positive rate:  $\frac{TP}{TP+FN} = \frac{TP}{n_+}$  (=Recall)

- False positive rate:  $\frac{FP}{TN+FP} = \frac{FP}{n_-}$

- Precision: How many selected items are relevant

$$\cdot \frac{TP}{TP+FP}$$

- Recall: How many relevant items are selected

$$\cdot \frac{TP}{TP+FN}$$

## Generative vs. Discriminative

- Used for classification

- Discriminative: Focus on the decision boundary

$$\cdot P(y|x)$$

· Sagt nicht über outliers aus

- Generative: What kind of positives/negatives are you likely to see

$$\cdot P(x,y) = p(y|x) f(x)$$

↳ How probable is it  
that  $x$  appears ( $\rightarrow$  Prior)

## Maximum Likelihood Estimation

- Point estimate, no representation of uncertainty  $\Rightarrow$  MAP
- $\theta_{MLE} = \operatorname{argmax} P(D|\theta)$
- Least squares regression

## Maximum a Posteriori

- Assumes a joint distribution  $P(D, \theta)$
- $P(x_1, \dots, x_n | \theta)$
- $\theta_{MAP} = \operatorname{argmax} P(\theta | D) = \operatorname{argmax} P(D | \theta) \cdot \underbrace{p(\theta)}_{\text{prior distribution}}$
- If prior distribution is uniform: MAP = MLE
- Allows to model prior belief about parameters
- Ridge Regression

## Gaussian Naïve Bayes

- Classification
  - $P(x, y) = p(x|y) \cdot P(y) = p(x_1|y) \cdots p(x_d|y) \cdot P(y)$
  - For new  $x$ , predict its  $y$
1. Compute probabilities
  2. Compute  $y = \operatorname{argmax} p(y|x) = p(y | \underbrace{x_1, \dots, x_n}_{\text{components of } x})$
- $$= P(y) \prod P(x_i | y)$$

$$P\left(\frac{\text{Banana}}{\text{Long, Sweet, Yellow}}\right) = \frac{P\left(\frac{\text{Long}}{\text{Banana}}\right) \times P\left(\frac{\text{Sweet}}{\text{Banana}}\right) \times P\left(\frac{\text{Yellow}}{\text{Banana}}\right) \times P(\text{Banana})}{P(\text{Long}) P(\text{Sweet}) P(\text{Yellow})}$$

} Kann weggelassen werden, da es

für alle  $y$   
den gleichen Wert hat.

- If model assumptions are met,  
 $\hookrightarrow$  NB will make same predictions

as Logistic Regression

## Gaussian Bayes

- we don't assume conditional independence

## Fisher's linear discriminant analysis (=generative)

- Models distribution of predictors separately in each of the response classes
- LDA is like PCA, but it focuses on maximizing the separability among known categories
- LDA creates a new axes where it can projects the data onto this axis in a way to maximize the separation of the two categories
- Maximize distance between the means of the categories
- Minimize scatter (=Variation) of categories

## Gaussian Mixture Model (GMM) (Unlabeled clustering)

- K-means: Only mean (= circles)
- Gaussian Bayes classifier: Trained with labels
- GMM:
  - each cluster is modeled as gaussian distribution with mean and variance (=oval)
  - we can say how probable a data is
  - Each point comes from a cluster with gaussian distribution
- we want to discover the variance/mean of each cluster with EM
  - ↳ Do MLE for each cluster to get parameters
- But we don't know which data-point goes where
- Hard EM:
  - clusters do not overlap
- Soft EM:
  - clusters may overlap

1. Place gaussians randomly in space
2. For each point figure out to which cluster it is most probable
  - ↳ Don't do hard assignment like k-means
3. Assign it to all cluster with probabilities
4. Recompute clusters
  - for partially labeled data:
    - for labeled data: assign directly to corresponding cluster
    - unlabeled data: Assign to most likely cluster

## Generative Adversarial Networks

- Unsupervised neural network learning
- Two neural nets run against each other
- Discriminator:
  - Has to decide whether the input data is from the real data set or generated from the generator
- Generator:
  - Generates data from random input noise
  - Learns the distribution of the real data-set
- Like a game between D and G

