

# IML Summary

Lasse Fierz - Ifierz

Version: 8. Februar 2023

## Basics

- General p-norm:  $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$
- Taylor:  $f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \mathcal{O}(x^3)$
- Power series of exp.:  $\exp(x) := \sum_{k=0}^{\infty} \frac{x^k}{k!}$
- $\sum_{k=0}^{\infty} (xy)^k = \frac{1}{1-xy}$
- Entropy:  $H(X) = \mathbb{E}_X [-\log \mathbb{P}(X=x)]$
- KL-Divergence:  
 $D_{KL}(P||Q) = \sum_{x \in \mathbb{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right) \geq 0$
- $1-z \leq \exp(-z)$
- Cauchy-Schwarz:  $|\mathbb{E}[X, Y]|^2 \leq \mathbb{E}(X^2)\mathbb{E}(Y^2)$
- Jensens Inequality: for a convex  $f(X)$ :  
 $f(\mathbb{E}(X)) \leq \mathbb{E}(f(X))$
- M p.s.d. if  $v^T M v \geq 0$

## Probability Theory:

- Gaussian:  $\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2})$
- $(N)(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu))$
- $X \sim \mathcal{N}(\mu, \Sigma), Y = A + BX \Rightarrow Y \sim \mathcal{N}(A + B\mu, B\Sigma^{-1}B^T)$
- Binomial Distr.:  $f(k, j; p) = \mathbb{P}(X = x) = \binom{n}{k} p^k (1-p)^{n-k}$
- $\mathbb{V}(X) = \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}(X^2) - [\mathbb{E}(X)]^2$
- $\mathbb{V}[X + Y] = \mathbb{V}[X] + \mathbb{V}[Y] + 2Cov(X, Y)$
- $Cov(X, Y) = \mathbb{E}[(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))]$
- $Cov(aX, bY) = abCov(X, Y)$

## Calculus

- $\int uv' dx = uv - \int u'v dx$  •  $\frac{\partial}{\partial x} \frac{g}{h} = \frac{g'h}{h^2} - \frac{gh'}{h^2}$
- $\frac{\partial}{\partial x} (b^T A x) = A^T b$  •  $\frac{\partial}{\partial x} (b^T x) = \frac{\partial}{\partial x} (x^T b) = b$
- $\frac{\partial}{\partial X} (c^T X^T b) = bc^T$  •  $\frac{\partial}{\partial X} (c^T X b) = cb^T$
- $\frac{\partial}{\partial x} (x^T A x) = (A^T + A)x \stackrel{\text{sym.}}{=} 2Ax$

- $\frac{\partial}{\partial X} Tr(X^T A) = A$  • Tr.trick:  $x^T A x \stackrel{\text{inner prod.}}{=} Tr(x x^T A) = Tr(A x x^T)$   
 $Tr(x x^T A) \stackrel{\text{cyclic perm.}}{=} Tr(A x x^T)$
- $\frac{\partial A^{-1}}{\partial x} = -A^{-1} \frac{\partial A}{\partial x} A^{-1}$
- $|X^{-1}| = |X|^{-1}$  •  $\frac{\partial}{\partial X} \log|x| = x^{-T}$  •  $\frac{\partial}{\partial x} |x| = \frac{x}{|x|}$
- $\frac{\partial}{\partial x} \|x\|_2 = \frac{\partial}{\partial x} (x^T x) = 2x$
- $\frac{\partial}{\partial x} \|x - b\|_2 = \frac{x-b}{\|x-b\|_2}$
- $\frac{\partial}{\partial x} \|x\|_1 = \text{sgn}(x)$
- $\sigma_{\text{sigmoid}}(x) = \frac{1}{1+\exp(-x)} \Rightarrow$
- $\nabla \sigma_{\text{sigmoid}}(x) = \sigma(x)(1 - \sigma(x)) = \sigma(x)\sigma(-x)$
- Jacobian =  $\frac{d}{dx} f(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}$
- Hessian:  $\nabla^2 f(x)$  e.g. for reg.  $2(X^T X + \lambda \mathbb{I}_d)$
- $\tanh x = \frac{2 \sinh x}{2 \cosh x} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- $\nabla \tanh x = 1 - \tanh^2 x$
- $\sin(a \pm b) = \sin(a)\cos(b) \pm \cos(a)\sin(b)$
- $\cos(a \pm b) = \cos(a)\cos(b) \mp \sin(a)\sin(b)$

## (Linear) Regression

General Regression: find  $\hat{y} = f(x) \leftrightarrow \min_{\hat{y}(x)} \|y - \hat{y}(x)\|_2^2$ .

$f^*(x)$  that minimizes L is  $\mathbb{E}[X|X=x]$  called **Bayes Optimal Predictor**. In practice unattainable.

Linear Regression: Weights are applied linearly:

$f(x) = \omega x$  or nonlinear **base fct**:  $f(x) = \omega \phi(x)$

Multidim.:  $L = \min_{\omega} \|Y - X\omega\|_2^2$ ,

$Y \in \mathbb{R}^n, x \in \mathbb{R}^{n \times d}, \omega \in \mathbb{R}^d$

## Closed Solution

If  $X^T X$  is invertible ( $X^T X$  has full rank  $\Leftrightarrow \text{rank}(X) = \min(d, n) \Rightarrow$  closed solution:  $\omega = (X^T X)^{-1} X^T Y$   
 $\nabla L$  is  $\mathcal{O}(nd)$ , closed solution is  $\mathcal{O}(nd^2)$ .

Can't apply closed solution for linearly dependent features.

**Note:** the closed solution can also be seen as finding the geom. proj. of  $y$  onto the hyperplane  $\text{span}(X)$ .  
 $(y - X\hat{\omega})^T X \omega = 0$

## Optimization

If not solvable in closed form or expensive to invert  $X^T X \rightarrow$  Gradient Descent:

$\omega_{t+1} \leftarrow \omega_t - \eta \nabla L(\omega_t)$ ,  $\eta$  is the learning rate.

Convergence guaranteed for  $\eta \geq \frac{2}{\lambda_{max}}$ , where  $\lambda_{max}$  is the

max EV of  $X^T X$ .

$X^T X$  diagonal  $\Rightarrow$  contour lines ( $L$  const) are ellipses

## Nonlinear Regression

Use fixed nonlinear feature maps of the inputs  $\phi(x)$  but still tune  $\omega \leftrightarrow \min \|y - \phi(x)\omega\|_2^2$ , with  $\phi(x) \in \mathbb{R}^{n \times p}$

**Note:** When working with NNNs both the weights and the nonlinear functions are chosen.

For closed solution same applies  $\text{rank}(\phi(x)) \stackrel{!}{=} \min(n, p)$

## Regularization

Among all unbiased solutions  $(X^T X)^{-1} X^T Y$  is the solution that has the smallest variance  $\Rightarrow$  minimizes gen. Error  
However the variance can get big  $\Rightarrow$  small  $L_{train}(\omega)$  but large  $L_{gen}(\omega)$  due to overfitting. Noise increases weights and regularization counters that effect.  $\Rightarrow$  Regularization:  
One can set the  $\omega$  of higher order features manually to zero ( $\Leftrightarrow$  choose a less complex model) or

## Ridge Regression

$\min_{\omega} \|Y - X\omega\|_2^2 + \lambda \|\omega\|_2^2$

Always allows for closed solution and lets LS converge faster through better conditioned problem (EVs of Hessian  $X^T X$  change)

Equivalent to performing Bayesianism approach with  $p(\omega) = \mathcal{N}(\omega|0, \Lambda^{-1})$  or linearly  $p(\omega) = \mathcal{N}(\omega|0, 1)$

Weights are decreased in general but not necessarily to exactly 0.

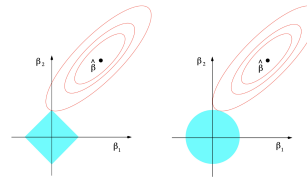
## Lasso Regression

Not a convex loss  $\Rightarrow$  no closed form solution

$\min_{\omega} \|Y - X\omega\|_2^2 + \lambda \|\omega\|_1$

Equivalent to performing Bayesianism approach with Laplacian prior:  $p(\omega_i) = \frac{\lambda}{4\sigma^2} \exp(-|\omega_i| \frac{\lambda}{\sigma^2})$

The weights of higher complexity features go to absolute zero  $\Rightarrow$  sparse weight vector result



Left: Lasso, Right: Ridge

In general with increasing  $\lambda$  the bias increases.  $\lambda_{opt}$  can be found using CV.

## Gradient Descent and Convexity

### Gradient Descent

$\omega_{t+1} \leftarrow \omega_t - \eta L(\omega_t)$

Converges to a stationary point.  $\nabla L(\omega) = 0 \Rightarrow$  GD stuck.

Complex fcts:  $\nabla L(\omega)$  from lin. approx. and use small  $\eta$

Large EVs for data depending heavily on one attribute and vice versa. Well conditioned if  $\lambda_{max}$  and  $\lambda_{min}$  are in similar range.

GD is sometimes slower and less accurate but there is more control and less comp. complexity

**Gradient Methods:** Momentum usage, Adaptive Methods, 2nd order methods

**Stochastic GD:** Use subsample from data for update step. Helps against saddle point conversion.

## Convexity

**Always:**

- global min/max  $\Rightarrow$  local min/max
- local min/max  $\Rightarrow$  stationary point

- $L(\omega) < L(v) \forall v \neq \omega \Leftrightarrow \omega$  is a global min

## Convexity:

- 0-order condition:  $L(sw + (1-s)v) \leq sL(w) + (1-s)L(v)$  aka function is lower or equal to linear connection of two points.
- 1st-order:  $L(v) \geq L(\omega) + \nabla L(\omega)^T (v - \omega)$  aka any point v on function is higher than point on linear approximation drawn at position  $\omega$
- 2nd-order:  $\nabla^2 L(\omega)$  is p.s.d. aka non-neg. curvature throughout function.
- $\omega$  stationary  $\Rightarrow \omega$  is local minimum
- $\omega$  is local minimum  $\Rightarrow \omega$  is global minimum

## Strong Convexity:

- 0-order:  $L(sw + (1-s)v) + \epsilon \leq sL(w) + (1-s)L(v)$  so fct always a bit below linear connection of points
- 1st-order: same as convex
- 2nd-order: strictly positive curvature always
- $\omega$  is global minimum  $\Rightarrow L(\omega) < L(v) \forall v \neq \omega$
- Only one global minimum

## Convexity Operations:

- Linear Comb. of convex functions are convex

- $f(g(x))$  is convex if f convex and g affine or f non-decreasing and g convex.

- Adding a convex and a strictly convex fct. yields a strictly convex function

## Model Selection

In general  $y = f(x) + \epsilon$ , where  $\epsilon$  is random noise

We can never know  $f(x)$  as we can only observe y. So we can't determine the estimation error  $(f(x) - \hat{f}(x))^2$

We use the gen. error  $(y - \hat{f}(x))^2 = \underbrace{(f(x) - \hat{f}(x))^2}_{\text{estimation error}} + \underbrace{\epsilon^2}_{\text{irreducible noise}} - \underbrace{2\epsilon(\hat{f}(x) - f(x))}_{0 \text{ on average}}$

Often interested in  $\mathbb{E}[(y - \hat{f}(x))^2] \approx \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$   
empirical error

## Bias and Variance

- Bias** =  $\mathbb{E}[(f(x) - \hat{f}(x))^2]$  Badness of model  
High for simple models and complex Ground Truths
- Variance** =  $\mathbb{E}[(\hat{f}(x) - \mathbb{E}\hat{f}(x))^2]$  fluctuation of  $\hat{f}$   
High for a too complex model and too little data (overfitting)

For the noiseless case  $y = f(x)$  a complex model can still overfit if the sample data is not representative of all data. Generalization Error = bias<sup>2</sup> + **variance**, idea of regularization: increase bias a bit to strongly decrease variance

#### Cross Validation

To estimate gen. error  $\Rightarrow$  train and test data. Usual splits are 50/50 and 80/20 (more often 80/20 because data is scarce)

To choose hyperparameters (e.g. regularization param  $\lambda$  or what choice of nonlinear features  $\phi(x)$ ) we perform k-fold cross validation: Split training data into k batches

- For each option of hyperparameter:
- for each batch:
  - Train model on the whole training data except for the batch
  - Calculate validation error on remaining batch
- Average validation error over all batches
- Choose hyperparameter with lowest avg. val. error
- Train model with that hyperparameter on the whole training set
- Determine test error

Leave one out CV (LOOCV):

- Split training data into sets of one  $\Rightarrow$  validation batch is of size 1
- Results in best model approximation
- Validation error is pretty bad (only one sample) but avg. ok
- Computationally expensive

#### Dataset Size

In general more data is always better. A limited dataset might not be representative of the underlying distribution. Usually  $y$  is noisy:  $y = f(x) + \epsilon$  in that case a small number of samples and a complex model will overfit the sample noise.

In the noiseless case  $n \rightarrow \infty \Rightarrow L_{train}(f(x)) \rightarrow 0$   
For  $n < d$  GD finds the solution that minimizes  $\|\omega\|_2$

#### Classification

- Probabilistic generative:  $p(x,y)$  allows for sample generation and outlier detection
- Prob. discriminative:  $p(y|x)$  classification with certainty
- Purely discr. c:  $X \rightarrow y$  just classification, easiest

Lin. seperable data  $\Rightarrow$  infinitely many solutions  $\Rightarrow$  SVM

#### Loss Functions

- Cross Entropy:  
$$\mathcal{L}^{CE} = - \left[ y' \log \hat{f}(x)' + (1 - y') \log (1 - \hat{f}(x)') \right]$$
Where  $y' = \frac{1+y}{2}$  and  $\hat{f}(x)' = \frac{1+\hat{f}(x)}{2}$

- Zero one loss:  $\mathbb{L}^{0/1} = \mathbb{I}\{\text{sign}(\hat{f}(x) \neq y)\}$   
Not convex nor continuous  $\Rightarrow$  surrogate logistic loss
- $\mathbb{L}^{\text{Hinge}} = \max(0, 1 - y\hat{f}(x))$
- $\mathbb{L}^{\text{percep}} = \max(0, -y\hat{f}(x))$
- $\mathbb{L}^{\text{logistic}} = \log(1 + \exp(-y\hat{f}(x)))$
- multidim. logistic loss: softmax:  
$$\mathbb{L}_i^{\text{softmax}} = \frac{e^{-a f_i}}{\sum_{j=1}^K e^{-a f_j}}$$
- $\mathbb{L}^{\text{exp}}(x)_i = \exp(-y\hat{f}(x))$

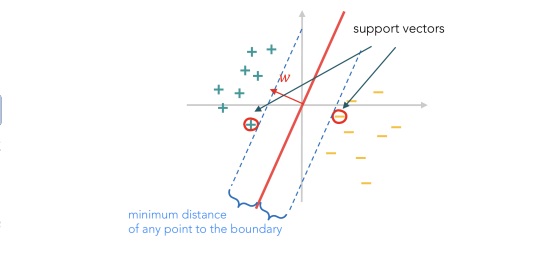
GD on logistic loss:  
$$\omega_{t+1} = \omega_t - \eta \frac{1}{n} \sum_{i=1}^n \nabla_{\omega} g(y \langle \omega_t, x \rangle) = \omega_t + \eta t \frac{1}{n} \sum_{i=1}^n \frac{y_i x_i}{1 + e^{y_i \langle \omega_t, x_i \rangle}}$$
 Converges to the  $\omega$  that minimizes the l2-distance to the decision boundary (SVM sol.)  
If classification error is not equally high for different classes  $\Rightarrow$  error metrics (see additional)

**Worst group error**(related to group fairness): Highest error among all clusters of a class (e.g. if one blob is 100% false)

**Robust generalization w.r.t. perturbations**  
Data augmentation, models that allow for invariance (e.g. CNNs)

**Distribution shifts** aka test data is different to training data: try to have the lowest possible error on the test samples that are similar to the training data.

#### SVM



Find  $\omega$  that maximizes the min distance of the closest points (support vectors) to the decision boundary.

margin =  $\min_i y_i \langle \omega, x_i \rangle$ , distance to SV =  $\frac{y_i \langle \omega, x_i \rangle}{\|\omega\|}$

Objective: maximize max margin direction:  $\arg\max_{\omega}$  margin( $\omega$ ) so that

Either  $\|\omega\| = 1$  or  $\|\omega\| = \frac{1}{\|\text{margin}\|}$   
Latter case: can look for  $\omega$  in the smaller subspace of  $\omega$  which yield a margin of 1

Objective:  $\mathcal{L}(\text{soft margin}) =$

$$\min_{\omega, \xi} \frac{1}{2} \|\omega\|^2 + C \sum_i \xi_i$$

s.t.  $y_i \omega^T x_i \geq 1 - \xi_i$  and  $\xi_i \geq 0 \forall i = 1, \dots, n$

Solve using lagrangian:

$$\mathcal{L} = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i \omega^T x_i)$$

#### Kernels

If we choose at least one nonlinear  $\phi(x)$  then  $\hat{f}(x)$  can be nonlinear

Note the comp. complexity of constructing  $\phi(x)$  (degree m polynomial of features  $X \in \mathbb{R}^{n \times d}$ ) is  $\mathcal{O}(nd^m) \Rightarrow$  huge for high dim. data

#### Kernel Trick

Feature maps only enter  $\hat{f}(x)$  by their inner product.  
Can write one of the possible global minimizers  $\hat{\omega} = \phi^T a$ ,  $a \in \mathbb{R}^n \Rightarrow$  Can write objective as:

$$L(\omega) = \frac{1}{n} \sum_{i=+}^n l$$

#### Other Nonlinear Models

##### K Nearest Neighbours

- For each datapoint determine the k nearest neighbours and assign a class based on the majority of the there present datapoints.
- Heavily dependent on k  $\Rightarrow$  Cross Validation
- Error prone in high dim. because of large distances  $\rightarrow$  optimize weights w.r.t. loss function (squared loss, cross-entropy etc.)
- Needs a lot of data but  $\mathcal{O}(nd)$  can be reduced to  $\mathcal{I}(n^p), p < 1$  if we allow for some error probability

##### Decision Trees

- At each node split data w.r.t. to one feature and a threshold (boundary at  $x_i > t_i$ )
- Each node returns class of the subset by majority
- Greedy Method: best step for current situation as opposed to generally best step  $\Rightarrow$  errors propagate.
- Very prone to overfitting as partitions can get very detailed
- $\Rightarrow$  random forest (averaged result over trees with random splits.)

#### Kernel Operations

Examples of valid kernels:

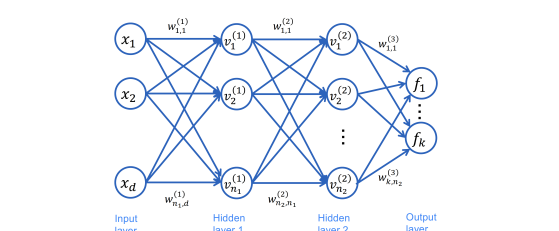
- $\alpha x^T x'$
- Polynomial:  $\alpha(x^T x' + \beta \mathbb{I})^p$
- RBF(Gaussian):  $\exp(-\frac{\|x - x'\|_2^2}{h^2})$
- Sigmoid:  $\tanh(\kappa x^T x' - b)$

Given two valid kernels  $k_1(x, x')$  and  $k_2(x, x')$  the following are also valid:

- $a k_1(x, x')$  for  $a \in \mathbb{R}$
- $k_1(x, x') + k_2(x, x')$
- $k_1(x, x') k_2(x, x')$
- $f(x) k_1(x, x') f(x')$
- $k(\phi(x), \phi(x'))$
- $g(k_1)$  with g: exp. or polyn. with all pos. coeff.

#### Neural Networks

NNs allow us to choose the feature maps in the model itself



$$z_1^{(1)} = \sum_{j=1}^d \omega_{1,j}^{(1)} x_j + \omega_{1,0}^{(1)}$$

$$z_k^{(l)} = \sum_{j=1}^{n_{l-1}} \omega_{k,j}^{(l)} v_j^{(l-1)} + \omega_{k,j}^{(l)} \quad v_l = \varphi(z_l)$$

for  $l = 1, \dots, L$  the number of layers  
Vector notation:

$$z^{(l)} = \mathbf{W}^{(l)} \mathbf{v}^{(l-1)} + \mathbf{W}_0^{(l)}$$

$$f(\mathbf{x}) = \mathbf{W}^{(L)} \mathbf{v}^{(L-1)}$$

where  $\mathbf{v}^{(l)} = [\varphi(z^{(l)}; 1)]$   $\varphi$  applied comp. wise

$\rightarrow$  optimize weights w.r.t. loss function (squared loss, cross-entropy etc.)

$$l(\mathbf{W}; \mathbf{x}, y) = \sum_{i=1}^n l_i(\mathbf{W}; \mathbf{x}_i, y_i)$$

Universal approx. theorem: Any cont. fct can be approximated by a finite layered NN with sigmoidal act. function.  
Weight decay reduces complexity

#### Activation Functions

A neural network with one hidden layer and nonlinear activation functions can approximate every continuous function.

- Sigmoid:  $\varphi(z) = \frac{1}{1 + \exp(-z)} = \frac{\exp(z)}{1 + \exp(z)}$ ,  
 $\varphi'(z) = \varphi(z)(1 - \varphi(z))$
- Relu:  $\varphi(z) = \max(0, z)$  (vanishing grad. for  $z < 0$ )
- Tanh:  $\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$
- ELU <sub>$\alpha$</sub>  (exp. relu):  $\begin{cases} \alpha(\exp(z) - 1), & \text{if } z < 0 \\ z, & \text{if } z \geq 0 \end{cases}$
- Softmax:  $\varphi(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

#### Backpropagation

Can reuse computations from **forward propagation** and from **layer  $l_{i+1}$  ... to compute  $W^{(i)}$**

$$(\nabla_{W^{(i)}} l)^T = \underbrace{\frac{\partial l}{\partial f}}_{\delta^{(L)}} \underbrace{\frac{\partial f}{\partial z^{(L-1)}}}_{\delta^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial z^{(L-2)}} \cdots \frac{\partial z^{(i+1)}}{\partial z^{(i)}} \underbrace{\frac{\partial z^{(i)}}{\partial W^{(i)}}}_{v^{(i-1)}}$$

$$\nabla_{W^{(l)}} l = \delta^{(l)} \mathbf{v}^{(l-1)}, \quad \delta^{(l)} = \text{diag}(\varphi'(z^{(l)}) W^{(l+1)}) \delta^{(l)}$$

$$\omega_{L-l}^{(l+1)} \leftarrow \omega_{L-l}^{(l)} - \eta \frac{\partial l}{\partial \omega_{L-l}^{(l)}}$$

**Note** usually minibatches are used for cum. weight updates.  
Running time grows linearly with num of params in feed forward

**Modifications:**



- $y = \operatorname{argmax}_{y'} P(y'|x) = \operatorname{argmax}_{y'} P(y') \prod_{i=1}^d P(x_i|y')$ 

$$(\pi_{1:k}^*, \mu_{1:k}^*, \Sigma_{1:k}^*) = \operatorname{argmin} - \underbrace{\sum_{i=1}^n \log \sum_{j=1}^k \pi_j \mathcal{N}(x; \mu_j, \Sigma_j)}_{\text{Intractable}}$$

With  $i$  the  $i$ th coordinate of a feature  $x$

- To predict:  $P(y|x) = \frac{P(y)P(x|y)}{\sum_y P(y)P(x|y)} \Rightarrow$  bayes. decis.
 

Can try SGD though but normally EM:
- Fisher's LDA:  $p = 0.5, \Sigma_+ = \Sigma_-$ .  $c = 2$ 

Expectation Maximization (EM)

 LDA if  $c \geq 2, p$  random, QDA = LDA with  $\Sigma_i \neq \Sigma_j$  Hard EM:

- Naive Bayes:  $\Sigma_y = \text{diag}(\sigma_{y,1}, \dots, \sigma_{y,d}), c \geq 2$
- GMMBC GMM instead of a single gaussian as likelihood estimation

**GNB:** Naive because it assumes independent samples (e.g. duplicate features would lead to overconfidence)

$$f(x) = \log \frac{p(Y=1|x)}{p(Y=-1|x)} \Rightarrow p(Y=1|x) = \frac{1}{1+\exp(-f(x))}$$

$$= \sigma(\omega^T x + \omega_0) \Rightarrow \text{same as log. reg. with}$$

$$\omega_0 = \log \frac{p_+}{1-p_+} + \sum_{i=1}^d \frac{\mu_{-,i}^2 - \mu_{+,i}^2}{2\sigma_i^2}, \omega_i = \frac{\mu_{+,i} - \mu_{-,i}}{\sigma_i^2}$$

**Note:** can introduce bias through choice of likelihood (e.g. GNB) to avoid overfitting  
 Also note that more clusters can fit the data better or as good (sup. no overfitting) while too few clusters is an issue.

Outlier Rejection

Can get  $p(x) = \sum_{y_i} p(x|y_i)p(y_i) \Rightarrow$   
 mark  $x_i$  as outlier if  $p(x_i) < \tau$

GANs

Discriminator working against Generator:

$$\min_{\omega_G} \max_{\omega_D} \underbrace{\mathbb{E}_{x \sim Data} \log D(x) + \mathbb{E}_{z \sim G} \log(1 - D(G(z)))}_{M(\omega_G, \omega_D)}$$

With  $D(x) = D(x, \omega_D), G(z) = G(z, \omega_G)$

- Finds saddle point
- For unlimited data  $D_G^* = \frac{P_{data}(x)}{P_G(x) + P_{data}(x)}$  Only conceptual as we can't know  $P_G$  and aren't given  $P_{data}$

Common training approach: simultaneous GD:

$$\omega_G^{(t+1)} = \omega_G^{(t)} - \eta_t \nabla_{\omega_G} M(\omega_G, \omega_D^{(t)})$$

$$\omega_D^{(t+1)} = \omega_D^{(t)} - \eta_t \nabla_{\omega_D} M(\omega_G^{(t)}, \omega_D)$$

Usually using minibatches of data

Possible problems: Oscillation, Mode collapsedata memo. leading to degeneracy, how to evaluate a GAN

Performance metric: Duality gap:

$$DG(\omega_G, \omega_D) := \max_{\omega_D} M(\omega_G, \omega_D') - \min_{\omega_G} M(\omega_G', \omega_D)$$

$$DG \geq 0, DG = 0 \text{ if } \omega_G \& \omega_D \text{ perform perfect equilibrium.}$$

If D and G have sufficient cap. DG upper bounds Jensen-Shanon Divergence

Gaussian Mixture Model

$$P(x|\theta) = P(x|\mu, \Sigma, \pi) = \sum_{j=1}^k \pi_j \mathcal{N}(x; \mu_j, \Sigma_j)$$

Note that GMMs are the same as GBC with unknown labels  
 MLE:

• PCA requires standardization because it considers the variance of the features in order to find the principle components.

Stdz always **after** train-test split.  
 Stdz not necessary for distance independent methods:

- Naive Bayes
- LDA
- Tree based methods (boosting, Random forests) etc.

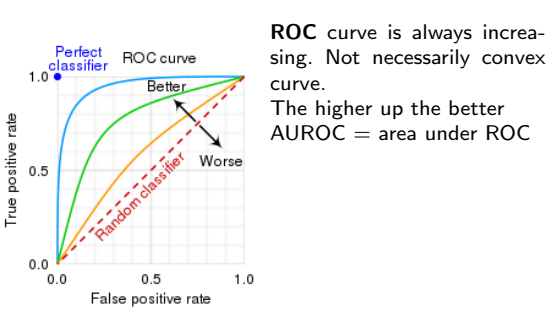
Classification Metrics

Define as positive the outcome which is crucial to get right.  
 Hypothesis test: Set hypothesis, reject it if  $\hat{p}(x) > \tau$  and accept it if  $\hat{p}(x) < \tau$   
 Reject hypothesis  $\Rightarrow$  positive — higher  $\tau \Rightarrow$  more negatives

- $\text{acc.} = \frac{TP+TN}{n}$
- $\text{prec.} = \frac{TP}{TP+FP}$
- $\text{FPR} = \frac{FP}{FP+TN}$
- $\text{Recall / TPR} = \frac{TP}{TP+FN}$
- $\text{balanced acc.} = \frac{1}{n} \sum_i TPR_i$
- $\text{FDR} = \frac{FP}{TP+FP}$
- $\text{F1-score} = \frac{2TP}{2TP+FP+FN}$
- $\text{ROC} = \frac{\text{FPR}}{\text{TPR}}$

$$\underbrace{\frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{prec.}}}}$$

**F1-score:** only high if both Recall and Precision are high  
 Useful if only interested in positive class  
 ROC curve:



**ROC** curve is always increasing. Not necessarily convex curve.  
 The higher up the better  
 AUROC = area under ROC

Individual Additions

Additional

Jensen's Inequality:  $\mathbb{E}_D \left[ \min_{f \in F} \hat{R}_D(f) \right] \leq \min_{f \in F} \mathbb{E}_D \left[ \hat{R}_D(f) \right]$

Standardization

Standardizing features  $x_{new} = \frac{x-\mu}{\sigma}$  yields values between 0 and 1. Necessary if one feature is much bigger than others and has a bigger influence on the weights. Standardizing allows for higher learning rates. Especially important for euclidian distance based methods like **knn,SVM,PCA,NNs,GD**

- KNN and SVM are methods based on the euclidian distance between the points
- NNs converge faster with standardized data. Also helps with vanishing gradients.