# IML Summary

Lasse Fierz - lfierz

Version: 8. Februar 2023

## Basics

- General p-norm: $||x||_p = (\sum_{i=1}^n |x_i|^p)^{1-p}$

- Taylor: $f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \mathbb{O}(x^3)$

- 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 $v^T M v \succeq 0$

Probability Theory:

- Gaussian: $\mathcal{N}(x|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}}exp(-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2})$

- $(N)(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d|\boldsymbol{\Sigma}|}}exp(-\frac{1}{2}(x-\mu)^T\boldsymbol{\Sigma}^{-1}(x-\mu))$

- $X \sim \mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma}), Y = A + BX \Rightarrow$
$Y \sim \mathcal{N}(A + B\boldsymbol{\mu}, B\boldsymbol{\Sigma}^{-1}B^T)$

- 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

- $\int uv'dx = uv - \int u'vdx$ $\quad \frac{\partial}{\partial x}\frac{g}{h} = \frac{g'h}{h^2} - \frac{gh'}{h^2}$

- $\frac{\partial}{\partial x}(\boldsymbol{b^T Ax}) = A^T b$ $\quad \frac{\partial}{\partial x}(\boldsymbol{b^T x}) = \frac{\partial}{\partial x}(\boldsymbol{x^T b}) = b$

- $\frac{\partial}{\partial \boldsymbol{X}}(\boldsymbol{c^T X^T b}) = \boldsymbol{bc}^T$ $\quad \frac{\partial}{\partial \boldsymbol{X}}(\boldsymbol{c^T X b}) = \boldsymbol{cb}^T$

- $\frac{\partial}{\partial x}(\boldsymbol{x^T Ax}) = (\boldsymbol{A^T + A})x \overset{\text{A sym.}}{=} 2\boldsymbol{Ax}$

- $\frac{\partial}{\partial \boldsymbol{X}}Tr(\boldsymbol{X^T A}) = A$ • Tr.trick: $\boldsymbol{x^T Ax} \overset{\text{inner prod.}}{=}$
$TR(\boldsymbol{x^T Ax}) \overset{\text{cyclic perm.}}{=} Tr(\boldsymbol{xx^T A}) = Tr(\boldsymbol{Axx^T})$

- $\frac{\partial \boldsymbol{A}^{-1}}{\partial x} = -\boldsymbol{A}^{-1}\frac{\partial \boldsymbol{A}}{\partial x}\boldsymbol{A}^{-1}$

- $|X^{-1}| = |X|^{-1}$ $\quad \frac{\partial}{\partial \boldsymbol{X}}log|x| = x^{-T}$ $\quad \frac{\partial}{\partial x}|x| = \frac{x}{|x|}$

- $\frac{\partial}{\partial \boldsymbol{x}}||x||_2 = \frac{\partial}{\partial \boldsymbol{x}}(\boldsymbol{x^T x}) = 2x$

- $\frac{\partial}{\partial \boldsymbol{x}}||\boldsymbol{x} - \boldsymbol{b}||_2 = \frac{\boldsymbol{x-b}}{||\boldsymbol{x-b}||_2}$

- $\frac{\partial}{\partial \boldsymbol{x}}||x||_1 = sgn(x)$

- $\sigma_{\text{sigmoid}}(x) = \frac{1}{1+exp(-x)} \Rightarrow$

- $\nabla\sigma_{\text{sigmoid}}(x) = \sigma(x)(1-\sigma(x)) = \sigma(x)\sigma(-x)$

- Jacobian $= \frac{d}{dx}f(\boldsymbol{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \cdots & & \cdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}$

- Hessian: $\nabla^2 f(x)$ e.g. for reg. $2(X^T X + \lambda\mathbb{I}_d)$

- $tanhx = \frac{2sinhx}{2coshx} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

- $\nabla tanhx = 1 - tanh^2 x$

- $sin(a \pm b) = sin(a)cos(b) \pm cos(a)sin(b)$

- $cos(a \pm b) = cos(a)cos(b) \mp sin(a)sin(b)$

## (Linear) Regression

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

$f^*(x)$ that minimizes L is $\mathbb{E}[X|\boldsymbol{X} = \boldsymbol{x}]$ 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.: $L = \underset{\omega}{min}||\boldsymbol{Y} - \boldsymbol{X\omega}||^2$,

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

### Closed Solution

If $X^T X$ is invertible ($X^T X$ has full rank $\Leftrightarrow rank(X) = min(d,n)) \Rightarrow$ closed solution: $\omega = (\boldsymbol{X^T X})^{-1}\boldsymbol{X^T Y}$
$\nabla L$ is $\mathbb{O}(nd)$, closed solution is $\mathbb{O}(nd^2)$.
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 span(X).
$(\boldsymbol{y} - \boldsymbol{X\hat{\omega}})^T X\omega = 0$

### Optimization

If not solvable in closed form or expensive to invert $X^T X \to$ Gradient Descent:
$\omega_{t+1} \leftarrow \omega_t - \eta\nabla L(\boldsymbol{\omega_t})$, $\eta$ 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 $\omega \leftrightarrow min||y - \phi(x)\omega||^2$, with $\phi(x) \in \mathbb{R}^{nxp}$

**Note:** When working with NNs both the weights and the nonlinear functions are chosen.
For closed solution same applies $rank\phi(x) \overset{!}{=} min(n,p)$

### 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**
$\underset{\omega}{min}||Y - X\omega||^2 + \lambda||\omega||^2$
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 $p(\omega) = \mathcal{N}(\omega|0,\boldsymbol{\Lambda}^{-1})$ or linearly $p(\omega) = \mathcal{N}(\omega|0,1)$
Weights are decreased in general but not necessarily to exactly 0.

**Lasso Regression**
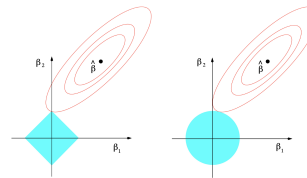Not a convex loss $\Rightarrow$ no closed form solution
$\underset{\omega}{min}||Y - X\omega||^2 + \lambda|\omega|$
Equivalent to performing Bayesianism approach with Laplacian prior: $p(\omega_i) = \frac{\lambda}{4\sigma^2}exp(-|\omega_i|\frac{\lambda}{w\sigma^2})$
The weights of higher complexity features 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 $(y - \hat{f}(x))^2 =$
$\underbrace{(f(x) - \hat{f}(x))^2}_{\text{estimation error}} + \underbrace{\epsilon^2}_{\text{irreducible noise}} - \underbrace{2\epsilon(\hat{f}(x - f(x)))}_{\text{0 on average}}$

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

### Bias and Variance

- **Bias** $= \mathbb{E}\left[(f(x) - \hat{f}(x))^2\right]$ Badness of model
High for simple models and complex Ground Truths

- **Variance** $= \mathbb{E}\left[(\hat{f}(x) - \mathbb{E}\hat{f}(x))^2\right]$ 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$^2$ $+$ **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

1. For each option of hyperparameter:

2. for each batch:

   - Train model on the whole training data except for the batch

   - Calculate validation error on remaining batch

3. Average validation error over all batches

4. Choose hyperparameter with lowest avg. val. error

5. Train model with that hyperparameter on the whole training set

6. 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 $n \to \infty \Rightarrow L_{train}(f(x)) \to 0$
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: $X \to y$
  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: $\mathbb{L}^{0/1} = \mathbb{I}\{sign(\hat{f}(x) \neq y)\}$
  Not convex nor continuous $\Rightarrow$ surrogate logistic loss

- $\mathbb{L}^{Hinge} = \max(0, 1 - y\hat{f}(x))$

- $\mathbb{L}^{percep} = \max(0, -y\hat{f}(x))$

- $\mathbb{L}^{logistic} = log(1 + exp(-y\hat{f}(x)))$

- multidim. logistic loss: softmax:
  $$\mathbb{L}_i^{softmax} = \frac{e^{-af_i}}{\sum_{j=1}^K e^{-af_j}}$$

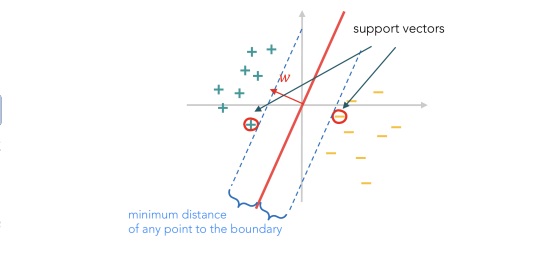- $\mathbb{L}^{exp}(x)_i = exp(-y\hat{f}(x))$

GD on logistic loss:
$\omega_{t+1} = \omega_t - \eta \frac{1}{n} \sum_{i=1}^n \nabla_\omega g(y\langle\omega_t, x\rangle) = \omega_t + \eta_t \frac{1}{n} \sum_{i=1}^n \frac{y_i x_i}{1+e^{y_i\omega_t^T x_i}}$ 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 additionals)
**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:
argmax margin$(\omega)$ so that
Either $||\omega|| = 1$ or $||\omega|| = \frac{1}{||margin||}$
Latter case: can look for $\omega$ in the smaller subspace of $\omega$ which yield a margin of 1
Objective: $\mathcal{L}$(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 $\xi_i \geq 0$ $\forall i = 1, .., n$
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 d features $X \in \mathbb{R}^{nxd}$)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

- Needs a lot of data but $\mathcal{O}(nd)$ can be reduced to $\ell(n^p), p < 1$ 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:
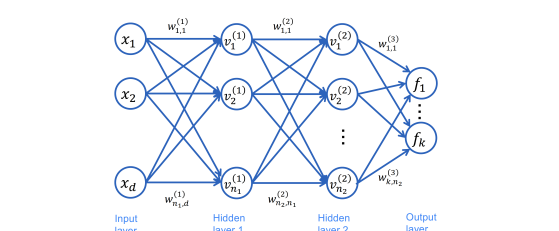
- $\alpha x^T x'$

- 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:

- $ak_1(x,x')$ 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



$$z_1^{(1)} = \sum_{j=1}^d \omega_{1,j}^{(1)} x_j + \omega_{1,0}^{(1)}$$
$$z_k^{(l)} = \sum_{j=1}^{n_{l-1}} \omega_{k,j}^{(l)} v_j^{(l-1)} + \omega_{k,j}^{(l)} \qquad v_l = \varphi(z_l)$$

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

$$\boldsymbol{z}^{(l)} = \boldsymbol{W}^{(l)}\boldsymbol{v}^{(l-1)} + \boldsymbol{W}_0^{(l)}$$
$$f(\boldsymbol{x}) = \boldsymbol{W}^{(L)}\boldsymbol{v}^{(L-1)}$$
where $\boldsymbol{v}^{(l)} = \left[\varphi(\boldsymbol{z}^{(l)}; 1)\right]$ $\varphi$ applied comp. wise

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

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

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$_\alpha$ (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 $W^{(i)}$

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

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

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

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

Momentum: initialize $d = 0$
$$\boldsymbol{d} \leftarrow m * \boldsymbol{d} + \eta_t \nabla_W l(\boldsymbol{W}; \boldsymbol{x}, y)$$
$$W \leftarrow W - \boldsymbol{d}$$

## Vanishing / Exploding Gradient

Potential reasons:

- $||\delta^{(i)}|| \rightarrow 0|\infty$ or $||v^{(i)}|| \rightarrow 0|\infty$

- Certain act. fct. like e.g. Relu (no saturation) can help avoid $\delta \rightarrow 0$

- Note $\delta$ only depends on $\varphi'$ while v depends on $\varphi$

- Helps to standardize input and / or use batch normalization

- Weight initialization matters as weight opt. is generally a non-convex problem

## Regularization in NNs

- Regularization term in Loss function

- Early stopping (before convergence to lowest training error)

- Dropout: deactivate about 50% of the nodes during training

- Data augmentation

**Batch normalization**
Normalize unit activations for a layer.
$BN(v, \gamma, \beta)$

- $\mu_s = \frac{1}{|S|}\sum_{i \in S} v_i$

- $\sigma_S^2 = \frac{1}{|S|}\sum i \in S(v_i - \mu_S)^2$

- $\hat{v}_i = \frac{v_i - \mu_S}{\sqrt{\sigma_S^2 + \epsilon}}$

- Scale and shift: $\bar{v}_i = \gamma\hat{v}_i + \beta$

Speeds up and stabilizes training. Solves covariate shift (different inputs also inbetween layers). Reduces importance of weight initialization. At initialization introduces exploding gradients. Has a mild regularization effect because of "random" batch size.

## Convolutional Neural Networks

Unfeasible to fully connect vectorized images due to number of parameters.

- Invariant regarding shifts, scale and rotation

- Updates still through backpropagation

- Still need nonlinear act. fct. for nonlinear functions

**Dimension of image after CNN layer** image: n x n, m kernels of size k*k, stride s, padding p:
$$l = \frac{n+2p-k}{s} + 1$$

## (Max)Pooling
Strongly reduces number of parameters

### Residual NNs

- Add possibility to skip layers e.g. feed input to intermediate layers

- Helps avoid vanishing gradients

- Allows for very deep NNs (1000+ layers)

- Can skip more than one layer (Dense Nets)

## Clustering

### K-Means

$$\text{Minimize } \hat{R}(\mu) = \sum_{i=1}^{n} \min_{j=1,..,k}||x_i - \mu_j||^2$$

**K-means algorithm / Lloyds heuristics:**
Initialize centers $\mu \rightarrow$ until convergence:
assign points to centers $\rightarrow$ relocate centers

- Risk function monotonically decreasing

- Converges to local minimum

- Cost per iteration: $\mathcal{O}(n * k * d)$
  With n points, k clusters and dim d

- Can use kernels for random shapes

- Strongly depends on initialization, Local convergence, How many clusters?

**K++ seeding**: Set one center $\mu_1 \rightarrow$ ad centers 2-k randomly with furthest datapoint from current clusters having highest prob. of becoming next center

$$\mu_j^{(0)} := x_i \text{ with prob. } \min_{l=1,..,j-1}||x_i - \mu_l^{(0)}||^2$$

Can be shown that expected loss is $(O)(logk)$ times that of opt. k-means solution
**Note:** finding optimal k is difficult. Can use heuristics
Generalization is very good(unsupervised learning)

## Dimensionality Reduction

Oftentimes we have high dimensional data which leads to high computational costs.
One countermeasure for this is dimensionality reduction.
$f: \mathbb{R}^D \rightarrow \mathbb{R}^d$, where $D > d$

### PCA

Finding linear transform to lower dimension that maximizes variance of data. Typically assume centered data.
$$\pi(x) = \mu^T x \Rightarrow \text{ maximize } \underbrace{\mu^T \frac{1}{n}\sum_i (\bar{x} - x_i)^2 \mu}_{Variance(\pi(x))}, ||\mu|| = 1$$

Note: we set $||\mu|| = 1$ to resolve uniqueness issue
$$\Rightarrow \mathcal{L} = \mu^T \frac{1}{n}\sum_i (\bar{x} - x_i)^2 \mu + \lambda(1 - ||\mu||)$$
$$\frac{\partial \mathcal{L}}{\partial \mu} \Rightarrow \underbrace{Var(x)}_{S}\mu = \lambda\mu \Rightarrow \lambda \text{ is highest EV of S}$$

$\mu$ is the respective Eigenvector
Note S becomes $\frac{1}{n}\sum_{i=1}^n x_i x_i^T$ when centered

Multiple dimensions: Solution can be optained in closed form also for $k > 1$

First PC with $X_1 = x$ as for $d = 1$, then
$$X_2 = X_1 - proj_{\mu_1} X_1$$
Then PCA with $X_2$, $X_3 = X_2 - proj_{\mu_2} X_2,...$
$$\Rightarrow \pi(x) = (x^T \mu_1, .., x^T \mu_d)$$

Alternatively: $\min_{W, ||W||^2=1} \sum_{i=1}^n ||x_i - z_i \boldsymbol{\omega}||^2 \rightarrow$ Regr.
$z_i^* = \boldsymbol{\omega}^T x_i \rightarrow \boldsymbol{\omega}^* = arg \min_{||\omega||_2=1} \sum_{i=1}^n ||\boldsymbol{\omega}\boldsymbol{\omega}^T x_i - x_i||_2^2$
$\rightarrow arg \min_{||\omega||_2=1} \boldsymbol{\omega}^T \boldsymbol{\Sigma}\boldsymbol{\omega}, \boldsymbol{\Sigma} = \sum_{i=1}^n \lambda_i v_i v_i^T$
PCA - K-means: Pretty much the same problem statement but for PCA $||\omega|| = 1$ and $z_i \in \mathbb{R}^k$ while for k-means $z_i$ are unit vectors

### Autoencoders

Aim at NN Identity encoding - decoding and then throw away the decoder.
Problem: Encoded features are not necessarily the best for ML task at hand. Features in latent space can take random form $\rightarrow$ bad for generation.
**Variational Auto Encoders:** AE with probabilistic latent space. E.g. yield parameters of distribution after encoding and sample from it for decoder. (also multiple latent reps. possible (HVAE))

## Statistical Perspective

- Quantify uncertainty

- make use of prior knowledge

- get access to other techniques

### Frequentism

Frequentist approach: $P(Y|X, \theta) \rightarrow MLE$
$\theta^* = \underset{\theta}{argmin} - log\sum_{i=1}^n p(y_i|x_i, \theta)$

- MLE is subject to consistency (param estimate converges to true params.)

- Asymptotic efficiency (smallest var $\forall$ well behaved estimators **for large n**)

- asymptotically normally distributed

Problem: might need a lot of data $\Rightarrow$ can overfit

### Bayesianism

Prior $p(\theta)$ about data, likelihood $p(y|x) \Rightarrow$ posterior $p(\theta|y, x) = \frac{p(\theta)p(y|x, \theta)}{p(y|x)}$, where p(x) are cancelled out.
Maximum Aposteriori Estimate (MAP):

$$\text{arg } \underset{\theta}{max} \; p(\theta) \prod_{i=1}^n p(y_i|x_i, \theta)$$
for $\theta \sim \mathcal{N}(0, \sigma^2 \mathcal{I})$ MAP yields ridge regr.

Can easily change priors e.g. laplacian prior $\Rightarrow$ lasso OLS
Gaussian yields very low prob. for values far from mean
laplacian: $p(x; \mu, b) = \frac{1}{2b}exp(-\frac{|x-\mu|}{b})$
student-t: allows more slack for furthe away values.
**General MAP procedure:**

- Choose likelihood fct. according to problem

- Choose prior for parameters

- Optimize MAP parameters

- Choose hyperparameters through CV

- Make predictions using Bayesian decision

### Bayesian Decision

Pick action that minimizes expected loss.

$$C(y, a) = \begin{cases} [y \neq a], & \text{if } a \in [+1, -1] \\ c, & \text{if } a = D \end{cases}$$

With abstention:

$$a^* = \begin{cases} y, & \text{if } P(y|x) \geq 1 - c \\ D, & \text{otherwise} \end{cases}$$

E.g. logistic loss in binary classification: Pick $y = +1$ if
$$p(y = -1|x) * \mathsf{C}_{FP} < p(y = +1|x) * \mathsf{C}_{FN} \Rightarrow$$
$$p(y = +1|x) \geq C_{FP}(C_{FP} + C_{FN})$$

Decision rule for bin. class.: (f is called Discriminative fct.)
$f(x) = log(\frac{p(y=1|x)}{p(y=-1|x)}) \Rightarrow f(x) > 0 \Rightarrow$ choose y=1
**Active Learning:** get labels in algorithm. Uncertainty sampling: Query labels of samples which the generative model is least sure about. i.i.d assumption is violated and can yield bad model. Can also query label of sample which when known makes other labels most probable.

### Conjugate Priors

Computing the posterior $p(\theta|x, y) = \frac{p(\theta)p(\theta|y,x)}{\int p(\hat{\theta})p(y|x, \hat{\theta})d\theta}$ Can be infeasible due to integration. $\Rightarrow$ conjugate priors make posterior of the same form as prior and int. is doable. (Almost no comp. cost)

| Prior/Posterior | Likelihood Fct |
|---|---|
| Beta | Bernoulli/Binom. |
| Dirichlet | Cathegorical/Multinom. |
| Gaussian, fixed cov. | Gaussian |
| Gaussian-inverse Wishart | Gaussian |
| Gaussian process | Gaussian |

Hyperparameters with CV

### Generative Modelling

Discriminative Modelling: $p(y|x)$, generative modelling: $p(y, x)$ strictly more ambitious goal
**Typical Approach:**

- Estimate prior on labels $p(y)$

- Estimate cond. distr. for each class $p(x|y)$ (e.g. how do the feature distributions look like)

- Joint distribution $p(y, x)$ through Bayes law

Alternative approach: $p(x) + p(x|y) \rightarrow p(x, y)$

### Gaussian Bayes Classifier

- Consider classes y with $p_y = p(Y = y), p_y = \frac{\#Y=y}{n}$

- $P(\boldsymbol{x}|y) \sim \mathcal{N}(\boldsymbol{x}; \mu_y, \sigma_y^2)$

$\boldsymbol{\mu}_y = \frac{1}{\#Y=y}\sum_{y_i=y}\boldsymbol{x}_i, \boldsymbol{\Sigma}_y = \sum_{y_i=y}(\boldsymbol{x}_i - \boldsymbol{\mu}_y)^2$
Aka $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ class wise

- $y = \underset{y'}{\operatorname{argmax}} P(y'|x) = \underset{y'}{\operatorname{argmax}} P(y') \prod_{i=1}^{d} P(x_i|y')$
  With i the ith 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.

- Fisher's LDA: $p = 0.5$, $\Sigma_+ = \Sigma_-$. $c = 2$
  LDA if $c \geq 2$, $p$ random, QDA = LDA with $\Sigma_i \neq \Sigma_j$

- Naive Bayes: $\Sigma_y = diag(\sigma_{y,1}, .., \sigma_{y,d})$, $c \geq 2$

- 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)
$f(x) = log\frac{p(Y=1|x)}{p(Y=-1|x)} \Rightarrow p(Y = 1|x) = \frac{1}{1+exp(-f(x))}$
$= \sigma(\boldsymbol{\omega}^T x + \omega_0) \Rightarrow$ 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:
$$\underset{\omega_G}{min}\ \underset{\omega_D}{max} \underbrace{\mathbb{E}_{x \sim Data} log D(\boldsymbol{x}) + \mathbb{E}_{\boldsymbol{z} \sim G} log(1 - D(G(\boldsymbol{z})))}_{M(\omega_G, \omega_D)}$$

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:
$$\boldsymbol{\omega}_G^{(t+1)} = \boldsymbol{\omega}_G^{(t)} - \eta_t \nabla_{\boldsymbol{\omega}_G} M(\boldsymbol{\omega}_G, \boldsymbol{\omega}_D^{(t)})$$
$$\boldsymbol{\omega}_D^{(t+1)} = \boldsymbol{\omega}_D^t - \eta_t \nabla_{\boldsymbol{\omega}_D} M(\boldsymbol{\omega}_G^{(t)}, \boldsymbol{\omega}_D)$$
Usually using minibatches of data
Possible problems: Oscillation, Mode collapsedata memo. leading to degeneracy, how to evaluate a GAN

Performance metric: Duality gap:
$$DG(\omega_G, \omega_D) := \underset{\omega_D'}{max} M(\omega_G, \omega_D') - \underset{\omega_G'}{min} M(\omega_G', \omega_D)$$
$DG \geq 0, DG = 0$ if $\omega_G \& \omega_D$ perform perfect equilibrium. If D and G have sufficient cap. DG upper bounds Jensen-Shanon Divergence

## Gaussian Mixture Model

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

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

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

Can try SGD though but normally EM:

## Expectation Maximization (EM)

Hard EM:

- Initialize parameters $\theta$

- E-step: Get most likely class for a datapoint
  $z_i^{(t)} = \underset{z}{\operatorname{argmax}} P(z|\boldsymbol{x}_i, \theta^{(t-1)})$
  After E-step we have complete datasets $\Rightarrow$ MLE

- M-step: MLE
  $\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} P(D^{(t)}|\theta)$

- Issues: Always works with more clusters ($\infty$ for many with $\sigma = 0$) also possible multiple on top of each other, doesn't consider certainties (bad for overlaps)

- k-means is like special case with uniform weights and spherical covariance

(soft) EM algorithm:

- Param initialization

- Estimate affiliation prob. of each datapoint and each cluster

- MLE

- Repeat

- Applications include: Clustering, Density estimation, Classification and Outlier detection

- Monotonically increases likelihood

- Depends on init. (often multiple random init. runs)

GMMs can overfit ($\infty$ clusters) can be fixed with prior or simpler model (Naive Bayes e.g.)

## Additionals

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

## 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.

- 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.
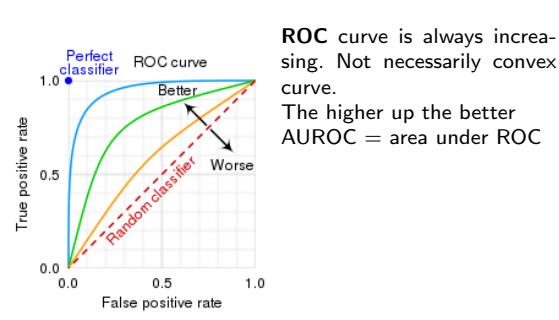
## Classification Metrics

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

- acc.$= \frac{TP+TN}{n}$    • prec.$= \frac{TP}{TP+FP}$

- FPR$= \frac{FP}{FP+TN}$    • Recall / TPR$= \frac{TP}{TP+FN}$

- balanced acc.$= \frac{1}{n}\sum_i TPR_i$   FDR$= \frac{FP}{TP+FP}$

- F1-score$= \underbrace{\frac{2TP}{2TP+FP+FN}}_{\frac{2}{\frac{1}{Recall}+\frac{1}{prec.}}}$   ROC$= \frac{FPR}{TPR}$

**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