

IML Summary

Lasse Fierz - Ifierz

Version: 9. 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 $v^T M v \succeq 0$

Probability Theory:

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

Calculus

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

- $\frac{\partial A^{-1}}{\partial x} = -A^{-1} \frac{\partial A}{\partial x} A^{-1}$
- $|X^{-1}| = |X|^{-1}$ • $\frac{\partial}{\partial X} \log|x| = x^{-T}$ • $\frac{\partial}{\partial x} |x| = \frac{x}{|x|}$
- $\frac{\partial}{\partial x} \|x\|_2 = \frac{\partial}{\partial x} (x^T x) = 2x$
- $\frac{\partial}{\partial x} \|x - b\|_2 = \frac{x-b}{\|x-b\|_2}$
- $\frac{\partial}{\partial x} \|x\|_1 = \text{sgn}(x)$
- $\sigma_{\text{sigmoid}}(x) = \frac{1}{1+\exp(-x)} \Rightarrow$
- $\nabla \sigma_{\text{sigmoid}}(x) = \sigma(x)(1 - \sigma(x)) = \sigma(x)\sigma(-x)$
- Jacobian = $\frac{d}{dx} f(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}$
- 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.

$f(x) = \omega x$ 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: $\omega = (X^T X)^{-1} X^T Y$

∇L is $\mathcal{O}(nd)$, closed solution is $\mathcal{O}(nd^2)$.

Alternatively geom. proj.: $(y - X\hat{\omega})^T X\omega = 0$

Optimization

If not solvable in closed form or expensive to invert $X^T X \rightarrow \omega_{t+1} \leftarrow \omega_t - \eta \nabla L(\omega_t)$, η is the learning rate.

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

max EV of $X^T X$.

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

Nonlinear Regression

Note: When working with NNs 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 One can set the ω of higher order features manually to zero (\Leftrightarrow choose a less complex model) or

Ridge Regression $\min \|Y - X\omega\|^2 + \lambda \|\omega\|^2$

Always closed solution and lets LS converge faster (EVs of Hessian $X^T X$ change)

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

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

Convex loss but no closed form solution (not differentiable)

Bayes. w/ Laplacian prior: $p(\omega_i) = \frac{\lambda}{4\sigma^2} \exp(-|\omega_i| \frac{\lambda}{w\sigma^2})$

$\omega_{\text{high}} \rightarrow 0 \Rightarrow$ sparse, λ_{opt} through 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 η

Large EVs \Leftrightarrow Dominant feature. Well conditioned if λ_{max}

and λ_{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: $L(sw + (1-s)v) \leq sL(\omega) + (1-s)L(v)$ 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 ω
- 2nd-order: $\nabla^2 L(\omega)$ is p.s.d. aka non-neg. curvature throughout function.
- ω stationary $\Rightarrow \omega$ is local minimum
- ω is local minimum $\Rightarrow \omega$ is global minimum

Strong Convexity:

- 0-order: $L(sw + (1-s)v) + \epsilon \leq sL(\omega) + (1-s)L(v)$ so fct always a bit below linear connection of points
- 1st-order: same as convex
- 2nd-order: strictly positive curvature always
- ω 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.
- Convex + strictly convex fct. = strictly convex fct

Model Selection

In general $y = f(x) + \epsilon$, where ϵ is random noise We can never know $f(x)$ as we can only observe y.

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

Bias and Variance

- Bias** = $\mathbb{E}[(f(x) - \hat{f}(x))^2]$
- Variance** = $\mathbb{E}[(\hat{f}(x) - \mathbb{E}\hat{f}(x))^2]$ fluctuation of \hat{f}

If $y = f(x)$ a complex model can still overfit. (Data sample)

Generalization Error = bias² + variance

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):
- Batch size 1
 - Results in best model approximation
 - Validation error bad (only one sample) but avg. ok
 - Computationally expensive

For $n < d$ GD finds solution that minimizes $\|\omega\|_2$

Classification

- Probabilistic generative: $p(x, y)$ (gen)
- Prob. discriminative: $p(y|x)$ (certainty)
- Purely discr. c: $X \rightarrow y$ (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:
$$\mathbb{L}_i^{\text{softmax}} = \frac{e^{-af_i}}{\sum_{j=1}^K e^{-af_j}}$$

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

GD on logistic loss: \rightarrow SVM sol.

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 augm., models that allow for invariance (e.g. CNNs)

Distribution shifts train \neq test data \Rightarrow try for sim. samples

SVM

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

argmax $\frac{\omega}{\|\omega\|}$ margin(ω).s.t. $\|\omega\| = 1$ or $\|\omega\| = \frac{1}{\|\text{margin}\|}$

Latter case: smaller subspace to look for ω

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

At least one nonlinear $\phi(x)$ then $\hat{f}(x)$ can be nonlinear

Complexity of constructing $\phi(x)$ is $\mathcal{O}(nd^m) \Rightarrow$ huge

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 assign majority class of knns

- 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 $\lambda(n^p), p < 1$ if we allow for some error probability

Decision Trees

- At each node split w.r.t. to one feature and threshold

- 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.

- Prone to overfitting as partitions can get very detailed

- \Rightarrow random forest (averaged over random splits)

Kernel Operations

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

For k1, k2 valid:

- $ak_1(x, x')$ for $a \in \mathbb{R}$
- $k(\phi(x), \phi(x'))$
- $k_1(x, x') + k_2(x, x')$
- $k_1(x, x')k_2(x, x')$
- $f(x)k_1(x, x')f(x')$
- $g(k_1)$ with g: exp. or polyn. with all pos. coeff.

Neural Networks

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

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

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

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

\rightarrow optimize weights w.r.t. loss fct (squared loss, CE 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 $_{\alpha}$ (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+1\$... to compute \$W^{\(i\)}\$](#)

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

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

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

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

$$\text{Momentum: initialize } d = 0$$

$$\mathbf{d} \leftarrow m * \mathbf{d} + \eta_t \nabla_{\mathbf{W}} l(\mathbf{W}; \mathbf{x}, y)$$

$$\mathbf{W} \leftarrow \mathbf{W} - \mathbf{d}$$

Vanishing / Exploding Gradient

Potential reasons:

- $\|\delta^{(i)}\| \rightarrow 0|\infty$ or $\|v^{(i)}\| \rightarrow 0|\infty$
- Relu (no saturation) can help avoid $\delta \rightarrow 0$
- Note δ only depends on φ' while v depends on φ
- Use standardized input and / or batch normalization
- Weight init. matters as weight opt. is gen. non-convex

Regularization in NNs

Reg. term in loss fct., early stopping, dropout or data augmentation

Batch norm. Normalize unit activations for a laye:

- $\mu_S = \frac{1}{|S|} \sum_{i \in S} v_i$
- $\sigma_S^2 = \frac{1}{|S|} \sum_i i \in S (v_i - \mu_S)^2$
- $\hat{v}_i = \frac{v_i - \mu_S}{\sqrt{\sigma_S^2 + \epsilon}}$
- Scale and shift: $\bar{v}_i = \gamma \hat{v}_i + \beta, \gamma, \beta$ given params.

Speeds up and stabilizes training. Solves covariate shift, reduces importance of weight init. Introduces exploding gradients. Mild reg. effect because of "random" batch size.

Convolutional Neural Networks

Less params than FC for images / audio

- 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

- Helps avoid vanishing gradients

- Allows for very deep NNs (1000+ layers)

Clustering

K-Means

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

K-means algorithm / Lloyd's 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: $\mathcal{O}(n * k * d)$, (k clusters)
- Can use kernels for random shapes
- **Depends on init, Local convergence, Choice of k?**

K++ seeding: Set one center $\mu_1 \rightarrow$, centers 2-k:

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

Expected loss is $\mathcal{O}(\log k)$ times that of opt. k-means sol.

Note: finding optimal k is difficult. Can use heuristics

Generalization is very good(unsupervised learning)

Dimensionality Reduction

PCA

Finding linear transform to lower dimension that maximizes variance of data. Typically assume centered data.

$\pi(x) = \mu^T x \Rightarrow$ maximize $\mu^T \frac{1}{n} \sum_i (\bar{x} - x_i)^2 \mu, \|\mu\| = 1$

Note: we set $\|\mu\| = 1$ to resolve uniqueness issue

$$\Rightarrow \mathcal{L} = \mu^T \frac{1}{n} \sum_i (\bar{x} - x_i)^2 \mu + \lambda(1 - \|\mu\|)$$

$$\frac{\partial \mathcal{L}}{\partial \mu} \Rightarrow \underbrace{\text{Var}(x)}_S \mu = \lambda \mu \Rightarrow \lambda \text{ is highest EV of } S$$

μ is the respective Eigenvector

Note S becomes $\frac{1}{n} \sum_{i=1}^n x_i x_i^T$ when centered

Multidim: Closed form solution also for $k > 1$

$$\text{First PC with } X_1 = x \text{ as for } d = 1, \text{ then}$$

$$X_2 = X_1 - \text{proj}_{\mu_1} X_1$$

Then PCA with $X_2, X_3 = X_2 - \text{proj}_{\mu_2} X_2, \dots$

$$\Rightarrow \pi(x) = (x^T \mu_1, \dots, x^T \mu_d)$$

Alternatively:

$$\min_{W, \|W\|^2=1} \sum_{i=1}^n \|\mathbf{x}_i - z_i \boldsymbol{\omega}\|^2 \rightarrow \text{Regr.}$$

$$z_i^* = \boldsymbol{\omega}^T \mathbf{x}_i \rightarrow \boldsymbol{\omega}^* = \arg \min_{\|\boldsymbol{\omega}\|_2=1} \sum_{i=1}^n \|\boldsymbol{\omega}^T \mathbf{x}_i - x_i\|_2^2$$

$$\rightarrow \arg \min_{\|\boldsymbol{\omega}\|_2=1} \boldsymbol{\omega}^T \boldsymbol{\Sigma} \boldsymbol{\omega}, \boldsymbol{\Sigma} = \sum_{i=1}^n \lambda_i v_i v_i^T$$

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

Nonlinear PCA: same with $\phi(x)$

Statistical Perspective

Frequentism

Frequentist approach: $P(Y|X, \theta) \rightarrow \text{MLE}$

$$\theta^* = \underset{\theta}{\text{argmin}} - \log \sum_{i=1}^n p(y_i | x_i, \theta)$$

- Consistency (Convergence 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):

$$\arg \max_{\theta} p(\theta) \prod_{i=1}^n p(y_i|x_i, \theta)$$

for $\theta \sim \mathcal{N}(0, \sigma^2 I)$ MAP yields ridge regr.

$$\text{laplacian: } p(x; \mu, b) = \frac{1}{2b} \exp(-\frac{|x-\mu|}{b})$$

- Choose prior & likelihood fct. according to problem

- Optimize MAP parameters

- Choose hyperparameters through CV

- Make predictions using Bayesian decision

Bayesian Decision

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

With abstention:

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

E.g. logistic loss in binary classification: Pick $y = +1$ if

$$p(y = -1|x) * C_{FP} < p(y = +1|x) * C_{FN} \Rightarrow \\ p(y = +1|x) \geq C_{FP}(C_{FP} + C_{FN})$$

Decision rule for bin. class.: (f is called Discriminative fct.)

$$f(x) = \log(\frac{p(y=1|x)}{p(y=-1|x)}) \Rightarrow f(x) > 0 \Rightarrow \text{choose } y=1$$

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 though

Conjugate Priors

Posterior $p(\theta|x, y) = \frac{p(\theta)p(y|x, \theta)}{\int p(\theta)p(y|x, \theta)d\theta}$ Can be infeasible \Rightarrow

Conjugate priors (Hyperparameters with CV):

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

Generative Modelling

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

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

$$\mu_y = \frac{1}{\#y} \sum_{y_i=y} x_i, \Sigma_y = \frac{1}{\#y} \sum_{y_i=y} (x_i - \mu_y)^2$$

Aka μ and Σ class wise

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

With i the ith coordinate of a feature x

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

- 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: $\Sigma_y = \text{diag}(\sigma_{y,1}, \dots, \sigma_{y,d}), c \geq 2$

- GMMBC GMM instead of a single gaussian

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:

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

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

Gaussian Mixture Model

$$P(x|\theta) = P(x|\mu, \Sigma, \pi) = \sum_{i=1}^k \pi_i \mathcal{N}(x; \mu_i, \Sigma_i)$$

GMMs are the same as GBC with unknown labels

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

Expectation Maximization (EM)

Hard EM:

- Initialize parameters θ

- E-step: Get most likely class for a datapoint

$$z_i^{(t)} = \operatorname{argmax}_z P(z|x_i, \theta^{(t-1)})$$

After E-step we have complete datasets \Rightarrow MLE

- M-step: MLE

$$\theta^{(t)} = \operatorname{argmax}_{\theta} P(D^{(t)}|\theta)$$

- ∞ clusters possible, many with $\sigma = 0$ or overlaps, doesn't consider certainties

- 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 inits runs)

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

Additional

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

Standardization

$x_{new} = \frac{x-\mu}{\sigma}$ - if one feature is bigger than others. Allows higher learning rates. Important for distance based methods: knn,SVM,PCA,NNs,GD Stdz always after train-test split.

Classification Metrics

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

$$\bullet \text{ acc.} = \frac{TP+TN}{n} \quad \bullet \text{ prec.} = \frac{TP}{TP+FP}$$

$$\bullet \text{ FPR} = \frac{FP}{FP+TN} \quad \bullet \text{ Recall / TPR} = \frac{TP}{TP+FN}$$

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

$$\bullet \text{ F1-score} = \frac{2TP}{2TP + FP + FN} \quad \text{ROC} = \frac{\text{FPR}}{\text{TPR}}$$

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

Individual Additions