

Lecture 2

$Q(\gamma, f(x)) \triangleq$ loss between label γ and estimate $f(x)$

Conditional expected risk (given X): $R(f, X) = \int Q(\gamma, f(x)) P(\gamma|x) d\gamma$

Total expected risk (over X, γ): $R(f) = \mathbb{E}_X [R(f, X)] = \int_X R(f, X) P(X) dX$
 $= \iint_{\gamma, X} Q(\gamma, f(x)) P(X, \gamma) dX d\gamma$

Empirical risk on training data $Z^{\text{train}} = \{(x_1, \gamma_1), \dots, (x_n, \gamma_n)\}$
 test data $Z^{\text{test}} = \{(x_{n+1}, \gamma_{n+1}), \dots, (x_{n+m}, \gamma_{n+m})\}$

Training error $\hat{R}(\hat{f}, Z^{\text{train}})$ for Empirical Risk Minimizer \hat{f} (ERM)

↳ select ERM $\hat{f} = \arg \min_{f \in \mathcal{C}} \hat{R}(f, Z^{\text{train}})$

↳ training error $\hat{R}(\hat{f}, Z^{\text{train}}) = \frac{1}{n} \sum_{i=1}^n Q(\gamma_i, \hat{f}(x_i))$

↳ test error $\hat{R}(\hat{f}, Z^{\text{test}}) = \frac{1}{m} \sum_{i=n+1}^{n+m} Q(\gamma_i, \hat{f}(x_i))$

$\hat{R}(\hat{f}, Z^{\text{test}}) \neq \mathbb{E}_X [R(\hat{f}, X)]$, as we can not guarantee that test set has the "real" distribution of X .

Taxonomy of Data

Measurement: Given an object set, a measurement X maps an object set into a domain \mathbb{K}

$$X: O^{(1)} \times \dots \times O^{(R)} \rightarrow \mathbb{K}$$
$$(o_1, \dots, o_R) \mapsto x_{o_1, \dots, o_R}$$

monadic data: $X: O \rightarrow \mathbb{R}^d$

d-adic data: $X: O^{(1)} \times O^{(2)} \rightarrow \mathbb{R}$, pairwise if $O^{(1)}$ is same space as $O^{(2)}$

polyadic data: $O^{(1)} \times O^{(2)} \times O^{(3)} \rightarrow \mathbb{R}$ + transformation invariances

Nominal scale

categories / stand on own

f is bijective

Ordinal scale

only meaningful in comparison with other samples e.g. rank

$f(x_1) < f(x_2)$ for $x_1 < x_2$

Interval scale

difference carries the information (Euclidean)

$aX + b$ at \mathbb{R}

Ratio scale

zero carries information, but not measurement unit (Newton)

aX at \mathbb{R}

Absolute scale

absolute values give meaning (γ) (logarithmic on axis)

Identity

Probability Spaces

Elementary event : $\omega_1 \dots \omega_N$ are sample points ω are coin flip

Sample space : $\Omega = \{\omega_1, \dots, \omega_N\}$ ~~set~~ Realization of coin flips

Family of sets : event A is a set of elementary events with

$$\hookrightarrow A \subset \Omega$$

$$\hookrightarrow \omega \in A \text{ or } \omega \notin A$$

Algebra of events: \mathcal{A} is algebra of events; for $A \subset \Omega$

$$\hookrightarrow \Omega \in \mathcal{A}$$

$$\hookrightarrow \text{if } A \in \mathcal{A} \wedge B \in \mathcal{A}$$

$$\text{then } A \cup B \in \mathcal{A} \wedge A \cap B \in \mathcal{A} \wedge A \setminus B \in \mathcal{A}$$

Probability of events: Assign weights $p(\omega_i)$ to all $\omega_i \in \Omega$ with

$$0 \leq p(\omega_i) \leq 1$$

$$\sum_i p(\omega_i) = 1$$

Probability of an event: $A \in \mathcal{A}$ with $P(A) = \sum_{\{\omega_i : \omega_i \in A\}} p(\omega_i)$

Probability model : (Ω, \mathcal{A}, P)

Lecture 3

Posterior

$$\text{Bayes Rule} : P(\text{model} | \text{data}) = \frac{P(\text{data} | \text{model}) P(\text{model})}{P(\text{data})} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}$$

Assume Model $\gamma = f(X, \theta) + \epsilon$, $\gamma \neq \text{output}$

↳ Bayes. view: X, Y, θ are random variables

↳ Parametric stat. st.: functional form of $P(X, Y | \theta)$ is given,
now estimate θ of $P(\text{data} | \text{model})$

↳ Non-parametric stat. st.: sample X, Y to estimate $P(\text{data} | \text{model})$

↳ Statistical learning theory: minimize empirical risk $\hat{R}(f, z^{\text{train}})$, don't estimate $P(\text{data} | \text{model})$

Bayesianism is based on model assumptions on data model

- evidence is huge to compute
- parameters are random variables
- output is a distribution, not a value

Frequentist: - maximal likelihood $\hat{\theta}$:

- Fisher information $\hat{I}(\theta^*)$ measure for information content of densities
- sampling theory
- hypothesis testing

Maximal likelihood method:

1. Define a parametric model
2. Define a likelihood as a funct. of the parametric model
3. Compute an estimator by maximizing the likelihood

ML is consistent: $\theta_{\text{ML}} \rightarrow \theta^*$ in probability as $n \rightarrow \infty$

asymptotically normal: $\frac{1}{\sqrt{n}}(\theta_{\text{ML}} - \theta^*)$ converges in distribution to a random variable with distribution $N(0, J^{-1}(\theta) I(\theta) J^{-1}(\theta))$.

asymptotically efficient: θ_{ML} minimizes $\mathbb{E}[(\theta_{\text{ML}} - \theta^*)^2]$ as $n \rightarrow \infty$

Rao-Cramér bound: $\mathbb{E}[(\hat{\theta} - \theta^*)^2] \geq \frac{1}{I(\theta^*)}$, for any unbiased estimator $\hat{\theta}$

↳ ML equals Rao-Cramér bound for $n \rightarrow \infty \Leftrightarrow$ asymptotically efficient

Bayesian methods

- allow priors
- provides a distribution of parameters when estimating
- requires efficient tricks when estimating posterior
- Prior often induces regularization

Frequentist methods

- No priors

- provides a single parameter when estimating

- requires only differentiation

+

intuition

$$Bias(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta^*$$

Left out: Rao-Blackwell derivation, consistency of ML
ML is also uniformly efficient (to θ^*)

Bayesian Learning

Posterior distribution: θ is considered as a random variable with $p(\theta|\mathcal{X})$

Assumption: $p(x)$ is unknown ($X \in p(x)$), $p(x|\theta)$ is known

"Wanted": $p(X=x|\mathcal{X}) = \int p(x|\theta) p(\theta|\mathcal{X}) d\theta$

Approximation: $p(\theta|\mathcal{X}) \sim \delta(\theta - \hat{\theta}) \Rightarrow p(x|\mathcal{X}) \approx p(x|\hat{\theta})$

Recursive Bayesian Estimator — also is an online algorithm

Assumption: Data are available in sequential order e.g. $\mathcal{X}^n = (x_1, x_2, \dots, x_n)$
are used one after another

Likelihood of data: $p(\mathcal{X}^n|\theta) = p(x_n|\theta) \underbrace{p(\mathcal{X}^{n-1}|\theta)}_{\text{similar}}$

posterior: $p(\theta|\mathcal{X}^n) = \frac{p(x_n|\theta) p(\theta|\mathcal{X}^{n-1})}{\int p(x_n|\theta) p(\theta|\mathcal{X}^{n-1}) d\theta} \underbrace{\text{similar}}$

prior: $p(\theta|\mathcal{X}^0) = p(\theta)$

Lecture 4

- (Parametric) maximum likelihood: Assume $y|X \sim N(f(X), \sigma^2 I)$

$$\text{solve: } \arg \max_{\beta} \sum_{i=1}^n \log P(Y=y_i | X=x_i, \beta)$$

- Statistical learning theory: $\arg \min_{f} \sum_{i=1}^n L(y_i - f(x_i))^2$ — same for linear regression

Linear regression model: $\hat{Y} = \beta_0 + \sum_{j=1}^d \beta_j x_j$ $X \in \mathbb{R}^d$ data, $Y \in \mathbb{R}$ output
 $\beta \in \mathbb{R}^d$ parameters, $\epsilon \sim N(0, \sigma^2)$

Residual sum of squares: $RSS(\beta) = (\hat{Y} - Y)^T (\hat{Y} - Y)$, $\hat{\beta} = (X^T X)^{-1} X^T Y$

Distribution of $\hat{\beta}$: $\hat{\beta} \sim N(\beta, (X^T X)^{-1} \sigma^2)$
 $\hat{\beta} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (X\beta + \epsilon)$, ϵ being the noise vector
 $= \beta + (X^T X)^{-1} X^T \epsilon$, only random variable, β is fixed

$$\mathbb{E}[\hat{\beta}] = \beta + (X^T X)^{-1} X^T \mathbb{E}[\epsilon] = \beta$$

$$\begin{aligned} \text{Cov}[\hat{\beta}] &= \mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T] \\ &= \mathbb{E}[(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}] \\ &= (X^T X)^{-1} X^T \mathbb{E}[\epsilon \epsilon^T] X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} \cancel{X^T X} \cancel{(X^T X)^{-1}} \end{aligned}$$

Optimality of Least Squares Estimates:

assume we want to find some $k = a^T \beta$, we only have $\hat{k} = a^T \hat{\beta}$

$$\mathbb{E}[a^T \hat{\beta}] = \mathbb{E}[a^T (X^T X)^{-1} X^T Y] = a^T (X^T X)^{-1} X^T (X\beta + \mathbb{E}[\epsilon]) = a^T \beta$$

a is unbiased

$$\text{Var}[a^T \hat{\beta}] = \text{Var}[a^T (X^T X)^{-1} X^T (X\beta + \epsilon)] = \text{Var}[a^T (X^T X)^{-1} X^T \epsilon]$$

$$= \mathbb{E}[a^T (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1} a] \stackrel{\text{no covariance}}{=} \sigma^2 a^T (X^T X)^{-1} a$$

We know define $q = a^T \gamma = a^T \hat{\beta} + a^T D \gamma = a^T ((X^T X)^{-1} X + D) \gamma = a^T \beta$
 Unbiased as seen above

$$\mathbb{E}[a^T \gamma] = a^T \mathbb{E}[\hat{\beta}] + a^T D \mathbb{E}[\gamma]$$

$$= a^T \beta + a^T D X \beta + a^T D \mathbb{E}[\gamma] = a^T \beta$$

$a^T \gamma$ unbiased

$$\hookrightarrow a^T D X = 0$$

Gauss-Markov Theorem - continuation of Optimality

For any linear estimator $\hat{\beta} = C^T \gamma$ that is unbiased for $\beta^T \beta$ we have

$$\text{Var}(C^T \beta) \leq \text{Var}(C^T \gamma)$$

$$\text{Var}(C^T \gamma) = C^T C \text{Var}(\gamma) C = C^T C^T S^2 I C$$

$$= S^2 C^T ((X^T X)^{-1} X^T + D) (X^T ((X^T X)^{-1} D^T) C)$$

$$= S^2 C^T ((X^T X)^{-1} X^T X ((X^T X)^{-1} D^T) C) \quad \begin{matrix} \text{cross-terms get eliminated} \\ \text{because } D X = 0 \end{matrix}$$

$$= S^2 C^T C S^2 ((X^T X)^{-1})^2 C + S^2 \|D\|_F^2 \|C\|^2$$

$$= C^T (C^T S^2 (X^T X)^{-1} C) + S^2 \|D\|_F^2 \|C\|^2$$

$$= C^T (\text{Var}(\beta)) C + S^2 \|D\|_F^2 \|C\|^2 \geq \text{Var}(C^T \beta) \quad \square$$

Claim: the least squares estimate $\hat{\beta}$ of β has the smallest variance among all linear unbiased estimators

Bias-variance trade-off $\hat{=}$ trade bias for variance reduction

Mean squared error = bias² + variance + noise variance

$$\mathbb{E}_{\gamma \mid X^*} [(\gamma - \hat{\gamma}(X^*))^2] = (\mathbb{E}[\gamma | X=X^*] - \mathbb{E}_\gamma [\hat{\gamma}(X^*)])^2 + \mathbb{E}[(\mathbb{E}_\gamma [\hat{\gamma}(X^*)] - \hat{\gamma}(X^*))^2]$$

$$+ \mathbb{E}_\gamma [\gamma - \mathbb{E}_\gamma [\gamma | X=X^*]]^2 |_{X=X^*}$$

Small data set, large set of functions \Rightarrow Variance large, bias small
Overfitting!

(good predictive)

Shrinkage / Regularization

Ridge regression

$$\min_{\beta} \sum_i (y_i - x_i^T \beta)^2 + \lambda \|\beta\|^2$$

+ protection against multicollinearity

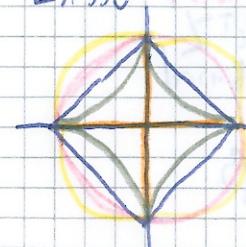
$$\text{Bayes: prior } \beta \sim N(0, S^2 I)$$

$$\text{Solution: } \hat{\beta}_R = (X^T X + \lambda I)^{-1} X^T y$$

mult. collinearity: $\hat{\beta}_R \hat{\beta} = \beta + (X^T X)^{-1} X^T \zeta$

LASSO

$$\text{Ridge: } \hat{\beta} = \beta + (X^T X + \lambda I)^{-1} X^T \zeta$$



can be combined to ElasticNet

Lasso

$$\min_{\beta} \sum_i (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_1$$

+ sparse \Rightarrow interpretable

$$\text{B.n.l. OVR } \beta_i = \frac{\lambda}{2S^2} \exp(-|\beta_i| \frac{\lambda}{2S^2})$$

Solution LARS

can be huge

$$\text{for } \hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^d x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^d |\beta_j|^q$$

Lecture 5

Bayesian:

$$(Y = X^T \beta + \epsilon \text{ with } \epsilon \sim N(\epsilon | 0, \delta^2)) \stackrel{\Delta}{=} N(Y | X^T \beta, \delta^2)$$

$$\text{Add R. dse Pr. or } p(\beta | \Delta) = N_d(\beta | 0, \Delta^{-1})$$

Model inversion: assume $p(\beta | X \gamma, \Delta) = N(\beta | \mu_\beta, \Sigma_\beta)$

$$\text{then } \mu_\beta = (X^T X + \delta^2 \Delta)^{-1} X^T Y$$

$$\Sigma_\beta = \delta^2 (X^T X + \delta^2 \Delta)^{-1}$$

Moments of Bayesian linear regression: of $Y = X\beta + \epsilon, \epsilon \sim N(\epsilon | 0, \delta^2 I_n)$

$$E_{\beta, \epsilon}[Y] = 0 \quad p(\beta | \Delta) = N_d(\beta | 0, \Delta^{-1})$$

$$\text{cov}[Y] = E_{\beta, \epsilon}[(X\beta + \epsilon)(X\beta + \epsilon)^T] = X\Delta^{-1}X^T + \delta^2 I_n$$

$$\text{cov}[Y] = X^{-1}(X X^T + \delta^2 I_n) \text{ if } \Delta = \lambda I_n$$

Moments of joint Gaussian: $E[Y] = 0, \text{cov}[Y] = X\Delta^{-1}X^T + \delta^2 I_n$

$$\left(\begin{array}{c} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{array} \right) \sim N \left(\begin{array}{c} 0 \\ \vdots \\ \vdots \\ 0 \end{array} \middle| \begin{array}{cccc} k_{1,1} + \delta^2 & k_{1,2} & \cdots & k_{1,n} \\ k_{2,1} & k_{2,2} + \delta^2 & \cdots & k_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ k_{n,1} & k_{n,2} & \cdots & k_{n,n} + \delta^2 \end{array} \right) \text{ with } k_{i,j} = k(x_i, x_j)$$

Kernelized linear regression $\stackrel{\Delta}{=} \text{Gaussian Process}$ $k(\cdot, \cdot) \stackrel{\Delta}{=} \text{kernel function}$

Good kernels specify the degree of similarity between any two data points.

More similar $\stackrel{\Delta}{=} \text{higher covariance} \stackrel{\Delta}{=} \text{less } \cancel{\text{axis aligned}} \text{ ellipse}$

Kernels

Properties: I Symmetry: $k(x_i, \cdot) = k(\cdot, x_i)$ feature's prediction

II Positive semi-definiteness: $\int k(x_i, \cdot) f(x) f(x) dx \geq 0 \quad \forall f \in \mathcal{H}$

Or Gram matrix $K = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix}$ is ~~pos.~~ positive semi-definite

Examples:

Linear	$x^T y$	$\begin{matrix} \text{linear} \\ \text{polynomial} \\ \text{Gauss/RBF} \\ \text{sigmoid/tanh} \end{math>$	$\begin{matrix} p \\ \exp(-\ x - y\ ^2/h^2) \\ \tanh(x^T y - b) \end{math}$	$\begin{matrix} \text{different} \\ \text{monomials} \\ \text{radial basis} \\ \text{functions} \end{math}$
poly	$(x^T y + 1)^p$ (or $p \in \mathbb{N}$)			
Gauss/RBF	$\exp(-\ x - y\ ^2/h^2)$			
Sigmoid/tanh	$\tanh(x^T y - b)$			

In kernel engineering, kernels can be combined in the following ways to result in another kernel

- addition $k_1(x_i, \cdot) + k_2(x_i, \cdot)$ also holds if kernels depend on x_j
- multiplication $k_1(x_i, \cdot) \circ k_2(x_i, \cdot)$
- scaling $k(x_i, \cdot) \cdot c$ with $c \geq 0$
- composition $f(k_1(x_i, \cdot))$ where f is a poly with positive coefficients or exp

Prediction 0-1 Gaussian processes

General: Cond. joint Gaussian distribution

$$\mathbb{P}(u_1, u_2) = N\left(\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \mid \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

$$\text{Then } \mathbb{P}(u_2 | u_1 = z) = N(u_2 | v_2 + \Sigma_{21} \Sigma_{11}^{-1} (z - v_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})$$

Gauss. process prediction: $\mathbb{P}(y_{n+1} | x_{n+1}, X_n) = N(y_{n+1} | \mu_{n+1}, \delta_{n+1}^2)$

$$\text{with } \mu_{n+1} = k^T (k + \delta^2 I)^{-1} y, \quad k = k(x_{n+1}, X)$$

$$\delta_{n+1}^2 = k(x_{n+1}, x_{n+1}) + \delta^2 - k^T (k + \delta^2 I)^{-1} k$$

Combining Regressors B estimators (simple average, $\hat{f}(x) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i(x)$)

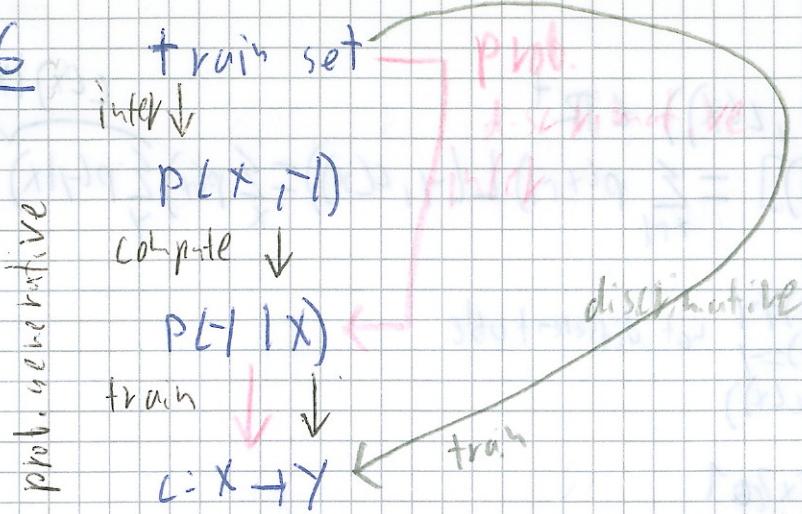
$$B \text{ avg: } \text{bias}[\hat{f}(x)] = \frac{1}{B} \sum_{i=1}^B \text{bias}[\hat{f}_i(x)] \Rightarrow \text{combined reduces variance}$$

$$\text{Variance: } \text{Var}[\hat{f}(x)] = \frac{1}{B^2} \sum_{i=1}^B \text{Var}_{\text{D}}[\hat{f}_i(x)] + \frac{1}{B^2} \sum_{i \neq j} \text{Cov}[\hat{f}_i(x), \hat{f}_j(x)]$$

\hookrightarrow If $\text{cov} = 0 \rightarrow$ Variance goes down by $\frac{1}{B}$

Good exercise

Lecture 6



Inter $p(x, l)$ \Rightarrow better understanding
new synthetic samples
outlier detection
More steps \Rightarrow more bias, and harder

$$\text{comp. } \log p(l | x) = 7 \text{ degrees of belief}$$

Inferring $p(x, l)$ I guess a family of parametrized probability models

Cause \downarrow Effect

$$\text{II} \text{ Inter } \theta \text{ with } \max_{\theta} \log p(X_{\text{train}} | l_{\text{train}}, x_{\text{train}} | \theta) - \max_{\theta} \log p(Y_{\text{train}} | \theta) \cdot p(X_{\text{train}} | Y_{\text{train}} | \theta)$$

(compute $p(l | x)$), with $p(x | l) = p(l)$ $p(x | l) = p(l) N(x | \mu_l, \Sigma_l)$

\hookrightarrow Linear discriminative models (LDA) iff $\Sigma_0 = \Sigma_1 = \Sigma_l$, classes have same size

$$p(l | x) = \frac{p(x | l) p(l)}{p(x | l) p(l) + p(x | 0) p(0)} \quad 0 \leq \text{class } 0, 1 \leq \text{class } 1$$

$$p(l | x) = \frac{1}{1 + \exp(-\log \frac{p(x | l) p(l)}{p(x | 0) p(0)})} = S(x^T w + b_0)$$

$$p(x | l) = (2\pi)^{-\frac{1}{2}} |\Sigma|^{-\frac{1}{2}} \exp(-\frac{1}{2} (x - \mu_l)^T \Sigma^{-1} (x - \mu_l)) \stackrel{\Delta}{=} \text{Normal distribution}$$

$$\log p(x | l) = -\frac{1}{2} \log (2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_l - \frac{1}{2} \mu_l^T \Sigma^{-1} \mu_l$$

$$\begin{aligned} \log \frac{p(x | l) p(l)}{p(x | 0) p(0)} &= \log(p(l)) - \log(p(0)) + \frac{\log(l)}{\log(0)} \\ &= x^T \underbrace{\Sigma^{-1} (\mu_l - \mu_0)}_w - \frac{1}{2} (\mu_l^T \Sigma^{-1} \mu_l - \mu_0^T \Sigma^{-1} \mu_0) + \frac{\log(p(l))}{\log(p(0))} \end{aligned}$$

\hookrightarrow generalized quadratic discriminant $\hat{\Sigma}_0 \neq \hat{\Sigma}_1 \Leftrightarrow \frac{1}{2} x^T \Sigma_1 x \neq \frac{1}{2} x^T \Sigma_0 x$

$$\log \frac{p(x | l) p(l)}{p(x | 0) p(0)} = x^T w x + w^T x + w_0 \Rightarrow p(l | x) = S(x^T w x + w^T x + w_0)$$

Bayes decision theory

I Define loss: $L(\gamma, c(x)) \in \mathbb{R}^+$

$$\text{II } \min_{\gamma} \mathbb{E}_{x \sim P} [L(\gamma, c(x))] = \sum_{x \sim P} p(x) L(\gamma, c(x)) = \sum_x p(x) \underbrace{\sum_y p(\gamma|y)}_{c(x)|x} L(\gamma, c(x))$$

Example losses

0-1 loss: $\begin{cases} 1 & \text{if } c(x) \neq y \\ 0 & \text{if } c(x) = y \end{cases}$ not differentiable

exponential loss: $\exp(-\gamma w^T x)$

Hinge: $\begin{cases} 0 & \text{if } \gamma w^T x \geq 1 \\ -1 - w^T x & \text{else} \end{cases}$

Perceptron loss: $\begin{cases} 0 & \text{if } \gamma w^T x \geq 0 \\ -1 - w^T x & \text{else} \end{cases}$

Maximum likelihood estimation best weights w w.r.t log-likelihood training set

$$\exp(L(w)) = p(\text{train} | w) = \prod_{i \leq n} p(y_i | x_i; w) = \prod_{i \leq n} p(y_i | x_i; w) p(x_i; w)$$

if training set \mathcal{D}

$$= \prod_{i \in \mathcal{D}} \frac{p(y_i)}{1-p(y_i)} \prod_{i \in \mathcal{D}} \delta(w^T x_i)^{y_i} (1-\delta(w^T x_i))^{1-y_i}$$

$$p(1|x)$$

$L(w)$ is analytically intractable, but differentiable \Rightarrow do gradient descent

Gradient descent

$$NL = -L$$

$$k=0$$

$$w^{(k)} = \text{random}$$

while $\|y(k) \nabla NL(w^{(k)})\| \geq \epsilon$:

$$w^{(k+1)} = w^{(k)} - \eta(k) \nabla NL(w^{(k)})$$

$$k=k+1$$

Newton's method $=$ ~~batch~~ ~~stochastic~~ updates no learning rule (polynomial)

$$w^{(k+1)} = w^{(k)} - \nabla^2 NL(w^{(k)})^{-1} \nabla NL(w^{(k)}) \quad \text{found by deriving Taylor on } w^{(k+1)}$$

$$\text{Good } y(k) = \frac{\|\nabla NL(w^{(k)})\|^2}{\nabla NL(w^{(k)})^\top \nabla^2 NL(w^{(k)}) \nabla NL(w^{(k)})}$$

derived using Taylor expansion

$$\text{and } y(k) = w^{(k+1)} - w^{(k)} \rightarrow \text{Exercise}$$

$\nabla^2 NL(w^{(k)})^{-1}$ is hard to compute

Lecture 7 & 8

Perceptron did not distinguish between class. fns, when all samples are correct
 \rightarrow rank

\hookrightarrow SVM wants max-margin

$i \leq n, j \leq m$

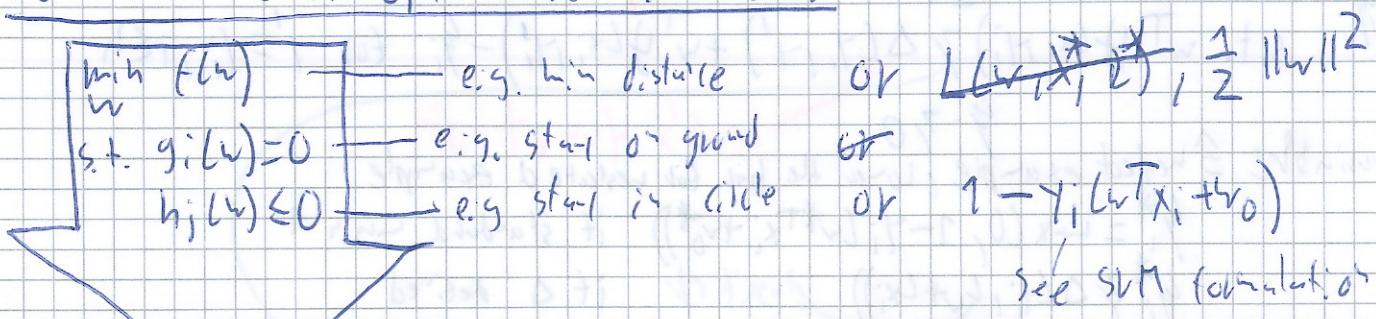
Slater's condition: $g_i(w) = 0, h_j(w) < 0 \Rightarrow$ Strong duality

Lagrangian: $L(w, \lambda, \alpha) = f(w) + \sum_i \lambda_i g_i(w) + \sum_j \alpha_j h_j(w)$, ≥ 0 SVM

complementary slackness: $\lambda_j^* h_j(w^*) = 0, w^* = \frac{1}{\lambda^*} \sum_i \lambda_i^* g_i(w^*)$

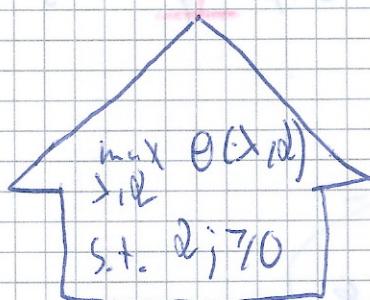
Extreme if strong duality (Slater's)

Constrained convex optimization via the dual



\rightarrow weak Duality / \rightarrow that strong if Slater's holds

$$\Theta(\lambda, \alpha) \geq \min_w L(w, \lambda, \alpha) \rightarrow \text{easy if with complementary slackness}$$



SVM - Maximum-margins classifier

x^+, x^- closest to border, proj. $\hat{\rightarrow}$ projection

$$\max_{w, w_0} 2 \ln(w^T w_0) = \|\text{proj}_w x^+ - \text{proj}_w x^-\| = \text{fixed}$$

$$= \left\| \frac{w^T x^+}{\|w\|^2} w - \frac{w^T x^-}{\|w\|^2} w \right\| = \frac{1}{\|w\|} \|w^T x^+ - w^T x^-\|.$$

$$\text{such that: } \gamma_i(w^T x_i + w_0) \geq 0 \text{ (correct class)} \quad \hat{\rightarrow} |w^T x^+ - w^T x^-|$$

Since both are invariant to scaling of w there exist w^* with
 $w^T x^+ + w_0^* = 1, w^T x^- + w_0^* = -1$ setting flat in gives the

$$\text{SVM formulation: } \min_w \frac{1}{2} \|w\|^2 \text{ s.t. } \gamma_i(w^T x_i + w_0) \geq 1$$

$$f(w) \quad 1 - \downarrow \hat{\rightarrow} \gamma_i(w)$$

Good, for bonds

SVM dual formulation $\max_{w, w_0} \min_{\gamma_i} L(w, w_0, \gamma)$ s.t. $\gamma_i \geq 0$

Start with $L(w, w_0, \gamma) = \frac{1}{2} \|w\|^2 + \sum_i \gamma_i (1 - y_i (w^T x_i + w_0))$, $\gamma_i \geq 0$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w^* = \sum_i \gamma_i x_i, \quad \frac{\partial L}{\partial w_0} = 0 \Rightarrow \sum_i \gamma_i = 0$$

Plugs into $\max_w \sum_i \gamma_i - \frac{1}{2} \sum_{ij} \gamma_i \gamma_j y_i y_j x_i^T x_j$, $\gamma_i \geq 0$ & $\sum_i \gamma_i = 0$

$$w_0^* = -\frac{1}{2} (w^* T x + w^* T x)$$

(can be omitted as often 0
to $y_i^T f(x_i) + b(f_i)$)

Soft-margin formulation

$$\min_w \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \gamma_i$$

$$\text{s.t. } w^T y(x_i, \gamma_i) \geq \Delta(\gamma_i, \gamma'_i) + w^T y(x_i, \gamma'_i) - \gamma_i \text{ for } \gamma'_i = -\gamma_i, i \leq n$$

slack variable $\hat{\gamma}_i \geq 0$ example, lower the bar for restricted example

$$\hat{\gamma}_i^* = \max(0, 1 - y_i (w^* T x_i + w_0^*)) \text{ if standard margin 1}$$

$$\hat{\gamma}_i^* \geq \Delta(\gamma_i, w^T x_i) \text{ if } \Delta \text{ defined}$$

Training algorithm $\epsilon \geq \text{tolerance}$

$$w=0$$

$$y=0$$

constraints = \emptyset

while constraints change:

for $i \leq n$

$$\gamma'_i = \arg \max_{\gamma'_i \neq \gamma_i} \{ \Delta(\gamma_i, \gamma'_i) + w^T y(x_i, \gamma'_i) \}$$

$$\text{if } w^T y(x_i, \gamma'_i) \geq \Delta(\gamma_i, \gamma'_i) + w^T y(x_i, \gamma'_i) - \gamma_i - \epsilon$$

$$\text{constraints} = \text{constraints} \cup \{ w^T y(x_i, \gamma'_i) \geq \gamma'_i \}$$

$$w, y = \text{solve} \left[\min_w \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_i \gamma_i \text{ s.t. constraints} \right]$$

One vs. rest

For each class compute compatibility function $f_k(x) = \arg \max_w f_k(x)$

- Lots of work, no inter-class boundary compared to normal SVM with one score for all classes

SVM with kernels

+ works with infinite-dimensional embeddings

- requires careful kernel selection

- adopts to structured classification

- requires feature engineering

+ formulated as quadratic programming which is efficient - with kernels, training always high scattered from curse of dimensionality