

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 $\hat{f} = \arg \min_{f \in \mathcal{F}_{\text{E}}}$ $\hat{R}(f, Z^{\text{train}})$

↳ 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 or expon.)

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 $y|X \sim N(f(X), \sigma^2 I)$

$$\text{solve: } \arg \min_{\beta} \sum_{i=1}^n \log P(Y=y_i | X=x_i, \beta^2),$$

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

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

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

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

$$\hat{\beta} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (X\beta + \epsilon)$$

$$= \beta + (X^T X)^{-1} X^T \epsilon$$

only random variable, β is fixed

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

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

Optimality of Least Squares Estimates:

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

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

\downarrow is unbiased

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

$$= \mathbb{E}[a^T (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1} a] = \sigma^2 a^T (X^T X)^{-1} a$$

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

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

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

$\downarrow a^T D X = 0$

Gauss-Markov Theorem - continuation of Optimality

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

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

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

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

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

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

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

Claim:

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

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

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

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

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

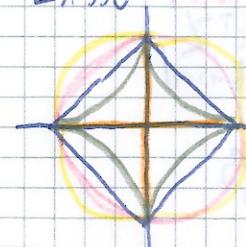
$$\text{Bayes: prior } \beta \sim N(0, S^2 I)$$

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

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

LASSO

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



$q=4$

$q=0.5$

$q=0.1$

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

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

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	$x^T y$	$\begin{matrix} \text{linear} \\ \text{polynomial} \\ \text{Gauss/RBF} \\ \text{sigmoid/tanh} \end{math>$	$\begin{matrix} p \\ \exp(-\ x - y\ ^2/h^2) \\ \tanh(x^T y - b) \end{matrix}$	$\begin{matrix} \text{different} \\ \text{monomials} \\ \text{radial basis functions} \\ \text{hyperbolic tangent} \end{math>$
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 $k(x_i, \cdot) \circ C$ 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, $\hat{f}(x) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i(x)$)

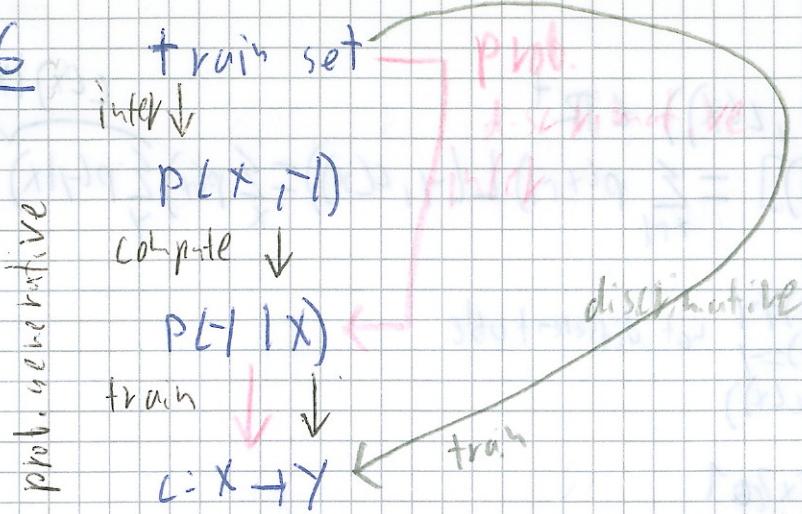
$$B \text{.us: } \text{bias}[\hat{f}(x)] = \frac{1}{B} \sum_{i=1}^B \text{bias}[\hat{f}_i(x)] \Rightarrow \text{combined regressors unbiased}$$

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

\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

$\left| \text{compute } p(l | x) = 7 \text{ degrees of belief} \right.$

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

Cause \downarrow
Effect

$$\text{II Infer } \theta \text{ with } \underset{\theta}{\max} \log p(X_{\text{train}} | X_{\text{train}}, l_{\text{train}} | \theta) = \underset{\theta}{\max} \log p(Y_{\text{train}} | \theta) \cdot p(X_{\text{train}} | Y_{\text{train}}, \theta)$$

(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

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

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

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

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

$$\begin{aligned} \log \frac{p(x | 1) p(1)}{p(x | 0) p(0)} &= \log(p(1)) - \log(p(0)) + \frac{\log(1)}{\log(0)} \\ &= x^T \underbrace{\Sigma_l^{-1} (\mu_1 - \mu_0)}_w - \frac{1}{2} (\mu_1^T \Sigma_l^{-1} \mu_1 - \mu_0^T \Sigma_l^{-1} \mu_0) + \frac{\log p(1)}{\log p(0)} \end{aligned}$$

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

$$\log \frac{p(x | 1) p(1)}{p(x | 0) p(0)} = x^T w x + w^T x + w_0 \Rightarrow p(l | x) = S(x^T w x + w^T x + w_0)$$

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: $\exp(-y c(x))$

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

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

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

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

$$k=k+1$$

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

$i \leq n, j \leq m$

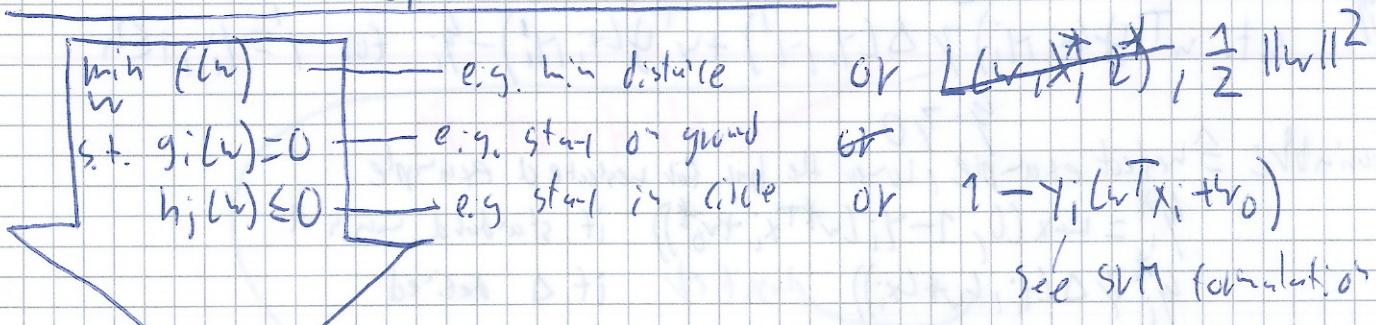
Slater's condition: $g_i(w) = 0, h_j(w) < 0 \Rightarrow$ Strong duality

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

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

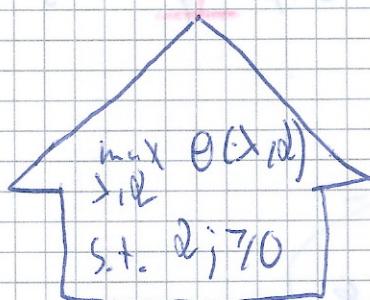
Extreme if strong duality (Slater's)

Constrained convex optimization via the dual



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

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



SVM - Maximum-margins classifier

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

$$\begin{aligned} \max_{w, w_0} 2 \ln(w^T w_0) &= \| \text{proj}_w x^+ - \text{proj}_w x^- \| = \text{dist} \\ \text{s.t. } w_i \geq 1, w_i \leq -1 &= \left\| \frac{w^T x^+}{\|w\|^2} w - \frac{w^T x^-}{\|w\|^2} w \right\| = \frac{1}{\|w\|} \|w^T x^+ - w^T x^-\|. \end{aligned}$$

such that: $y_i (w^T x_i + w_0) \geq 1 \text{ (correct class)} \quad \hat{\rightarrow} |w^T x^+ - w^T x^-|$

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

$$\begin{aligned} \text{SVM formulation: } \min_w \frac{1}{2} \|w\|^2 \text{ s.t. } y_i (w^T x_i + w_0) \geq 1 \\ f(w) \quad 1 - \downarrow \hat{\rightarrow} \hat{\geq} h_i(w) \end{aligned}$$

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 $\epsilon \geq \text{tolerance}$

$$w=0$$

$$y=0$$

$$\text{constraints} = \emptyset$$

while constraints change:

for $i \leq n$

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

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

$$\text{constraints} = \text{constraints} \cup \{ w^T y(x_i, \gamma_i) \geq \}$$

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

Prediction: $C(x) = \arg \max_y w^T y(x, y)$

One vs. rest

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

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

SVM with kernels

+ works with infinite-dimensional embeddings

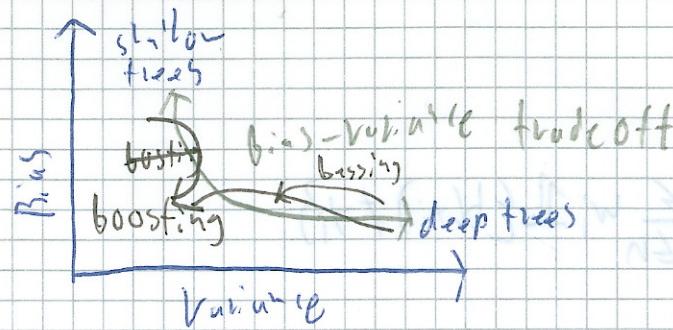
- requires careful kernel selection

- adopts to structured classification

- requires feature engineering

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

Lecture 9



Wise choice of words: If estimates are diverse, independent, decentralized, well aggregated and trusted, then the estimate aggregation will be better than that of each on its own.

Bagging: Bootstrap Aggregation:
 I Draw M bootstrap sets Z'_1, \dots, Z'_M
 II Train M base models $b^{(1)}, \dots, b^{(M)}$
 III Aggregate $\hat{b}^{(M)}(x) = \frac{1}{M} \sum b^{(i)}(x)$ regression
 majority, $(b^{(i)}(x))$ classification

Theorem: $\mathbb{E}_{\substack{x_1 \\ Z_1, Z_2, \dots, Z_M}} [(Y - \hat{b}^{(M)}(x))^2] \leq \mathbb{E}_{\substack{x_1, x_2 \\ Z_1, Z_2}} [(Y - b(x))^2]$

Proof

$$\begin{aligned} \mathbb{E}_{\substack{x_1, x_2 \\ Z_1, Z_2}} [(Y - b(x))^2] &= \mathbb{E}_{\substack{x_1, x_2 \\ Z_1, Z_2}} [(Y - \mathbb{E}_2[b(x)]) + \mathbb{E}_2[b(x)] - b(x))^2] \\ &= \mathbb{E}_{\substack{x_1 \\ Z_1}} [(Y - \mathbb{E}_2[b(x)])^2] + \mathbb{E}_x [(\mathbb{E}_2[b(x)] - b(x))^2] + \dots \\ &\approx \mathbb{E}_{\substack{x_1 \\ Z_1}} [(Y - b^{(M)}(x))^2] + \mathbb{E}_x [\text{Var}_2(b(x))] \\ &\geq \mathbb{E}[(Y - b^{(M)}(x))^2] \end{aligned}$$

for

$$\mathbb{E}[b^{(M)}(x)] = \mathbb{E}_2[b(x)] \quad \& \Pr_{Z_1, \dots, Z_M} (\|\mathbb{E}_2[b(x)] - b^{(M)}(x)\| < \epsilon) \xrightarrow{M \rightarrow \infty}$$

Each classifier is trained ~~independently on different~~ requires independence but the training sets are dependent, the correlation is small though.

① on different training sets which ensures diversity

Random Forests: bootstrap a set, but train a tree only on m features, not on all

AdaBoost $\triangleq \min_b \frac{1}{n} \sum L(t_i, b(x_i))$ with $\exp(-t_i) / \sum \exp(-t_i) = \alpha_i \cdot b^{(t)}(x)$

$$b^{(0)} = 0$$

$$w_1 = \frac{1}{n}$$

for t in $[1 \dots M]$

$$\text{train } b^{(t)} = \arg \min_b L^w(b) = \arg \min_b \sum_{i=1}^n w_i \mathbb{I}\{b(x_i) \neq t_i\}$$

$$\text{eval } e_{t+1} = L^w(b^{(t)})$$

$$\text{add } b^{(t)} = b^{(t-1)} + \alpha_t b^{(t)}, \text{ with } \alpha_t = \log\left(\frac{1}{e_{t+1}} - 1\right)$$

$$\text{reweight } w_i = w_i \exp(\alpha_t \mathbb{I}\{b^{(t)}(x_i) \neq t_i\})$$

normalize $w_1 \dots w_M$

end

Forward stagewise additive modeling (FSAM) \Leftarrow AdaBoost

$$\min_{b^{(M)}} \dots \min_{b^{(1)}} \frac{1}{n} \sum L(t_i, \alpha_0 b^{(0)}(x_i) + \alpha_1 b^{(1)}(x_i) + \dots + \alpha_M b^{(M)}(x_i))$$

AdaBoost trains max-margin classifiers

$$b^{(M)}(x_i) = \operatorname{sign}\left(\sum_{t=1}^M \frac{\alpha_t}{\sum_{t=1}^M \alpha_t} b^{(t)}(x_i)\right), \text{ margin} = t_i b^{(t)}(x_i) \xrightarrow{M \rightarrow \infty} 1$$

Boosting does not overfit

I with prob $\gamma \geq \frac{1}{2}$, gen. error $b^{(t)} \leq [\text{fraction of examples with margin} \leq \mu] + O\left(-\frac{1}{n} \log(H) - \frac{1}{\mu n}\right)$

II fraction of examples w/ margin $\geq \mu$

increases exponentially with H

shinks with grows until shinks with big
big n big H with big M

Random forests and AdaBoost truth

- spiky interpolating \leftarrow mostly smooth
sharp localized changes

- self-covering models to interpolate noisy examples

\Leftarrow trained model is an average of interpolating models

smooth flanks to closure of envelopes
interpolating \triangleq no training error

complex base \rightarrow spiky

\Leftarrow average further localizes noise interference

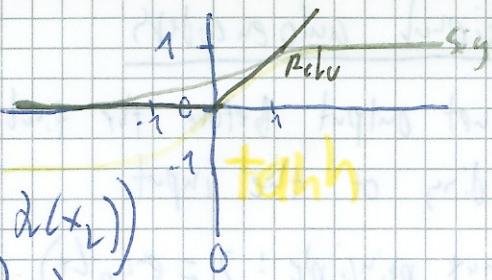
	short vertical stumps	stumps	tall trees	flattened NN not spiky
expected	high	medium	low	
v.s.	✓	✓	✓	
spiky	X	~	✓	

Lecture 10

Linear function $L(x) \mapsto w_x$

Activation function $d(x) \mapsto (d_1(x_1), \dots, d_n(x_n))$

Neural Net = $d^{(1)}_0 L^{(1)}_0 \circ \dots \circ d^{(n)}_0 L^{(n)}_0 L^{(n)}(x)$



Deep neural networks can approximate complex non-linear functions.

$f: [0,1]^n \rightarrow \mathbb{R}$ is continuous, then there is a NN such that

$$\text{I } \text{NN}(x) = \sum_{i \leq h} d_i \delta(w_i^T x + b_i) \quad \Rightarrow \text{dimensions of NN can be huge}$$

$$\text{II } \max_x |f(x) - \text{NN}(x)| < \epsilon$$

$$\text{softmax: } -P(y_i) = \frac{e^{p(\beta z_i)}}{\sum_{j \leq K} e^{p(\beta z_j)}}$$

(decision layer at end)

Training a NN: $\min_{\theta} \sum_{i \leq h} L(y_i, \text{NN}_{\theta}(x_i))$ → analytically intractable, but easy to differentiate

mini-batch (stochastic)

Gradient descent

repeat

do forward/backprop to compute

$$\frac{\partial \text{NN}}{\partial w_r^{(l)}} \Big|_{\theta, x_i} \quad \text{for each } l, r \text{ and } i \leq h \text{ and } i \in \text{sample set } S$$

$$\frac{\partial L}{\partial w_r^{(l)}} = \sum_{i \leq h} \frac{\partial L(y_i, \text{NN}_{\theta}(x_i))}{\partial w_r^{(l)}} \approx \sum_{i \in S} \frac{\partial L(y_i, \text{NN}_{\theta}(x_i))}{\partial w_r^{(l)}}$$

$$w_r^{(l)} = w_r^{(l)} - \eta(k) \frac{\partial L}{\partial w_r^{(l)}}$$

$k \leftarrow k + 1$

until test loss decreases significantly

If $\eta(k) \geq 0, \sum_k \eta(k) = \rho, \sum_k \eta^2(k) < \rho$ then gradient descent works

according to Polya's-Monte Carlo method

Full Batch Gradient Descent

+ better gradient precision

+ smaller generalisation error

Stochastic gradient descent

+ can handle large training sets

+ faster (in practice)

+ escapes local minima

Variational autoencoders

the net output before the final classification can be seen as a representative encoding of the input.

Information principle: $z = \text{enc}_\theta(x)$ is good if it maximizes $I(x, z) = \mathbb{E}_{x, z} [\log \frac{p(x, z)}{p(x)p(z)}]$

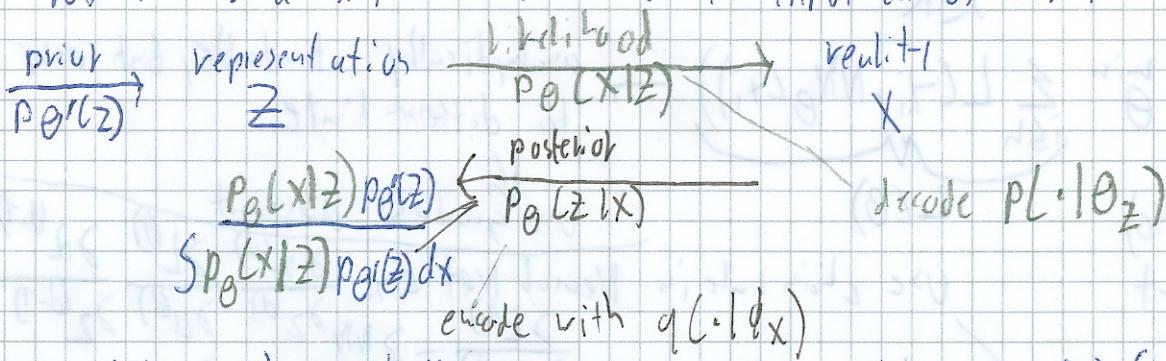
$$\text{argmax}_\theta I(x, z) = \mathbb{E}_{x, z} [\log p(x|z)] - \mathbb{E}_{x, z} [\log p(z)]$$

↓ gives each value in training set a unique representation, but meaningless

$\sum_i \mathbb{E}_{z|x_i} [\log p(z|x_i)]$ enc _{θ} (x_i)

Good representations are

- informative: you can guess the image given the representation
- disentangled: modifying one feature in the representation induces change in only one perceptible feature in the reconstruction
- robust: a slight modification in the input causes a slight change in the representation



Train VAE with argmax θ/θ_0 $\sum_{i \in n} \log p_{\theta/\theta_0}(x_i)$, where $p_{\theta/\theta_0}(x_i) = \int p_{\theta/\theta_0}(z) p_\theta(x|z) dz$

$$= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\log \left(\frac{p_{\theta/\theta_0}(x_i|z)}{p_{\theta/\theta_0}(z|x_i)} \right) \right] = \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\log \left(\frac{q_\phi(z|x_i)}{p_{\theta/\theta_0}(z|x_i)} \right) \right]$$

$$= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\log \left(\frac{p_{\theta/\theta_0}(x_i|z)}{q_\phi(z|x_i)} \right) \right] + \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\log \left(\frac{q_\phi(z|x_i)}{p_{\theta/\theta_0}(z|x_i)} \right) \right]$$

$\text{elbo}_{\theta/\theta_0}(\cdot)$ $\text{KL}(q_\phi(\cdot|x_i) \mid p_{\theta/\theta_0}(\cdot|x_i)) \geq 0$

$$\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] + \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log \left(\frac{p_\theta(z)}{q_\theta(z|x)} \right)]$$

Info term $-\text{KL}(q_\phi(\cdot|x_i) \mid p_{\theta/\theta_0}(\cdot)) \stackrel{\text{↑ regularization term}}{=} \text{disentanglement}$

estimates informativeness
robustness by choice of $p_\theta(\cdot|z)$ and $q_\phi(\cdot|x)$

Lecture 11

Multiple NN layers: some functions can not be efficiently represented
(by shallow architectures)

Consequences: we don't need exponentially many nodes per layer
- poor generalisation may be expected when using a too
deep shallow representation

Rule of thumb: 5-10 training examples per independent weight at least

Generative models goal: learn density $p_{data}(x)$ close to $p_{fit}(x)$

direct and generative models
solve for $p_{data}(x)$

learn submodel for $p_{data}(x)$
without explicitly defining it.

explicit densities

implicit densities

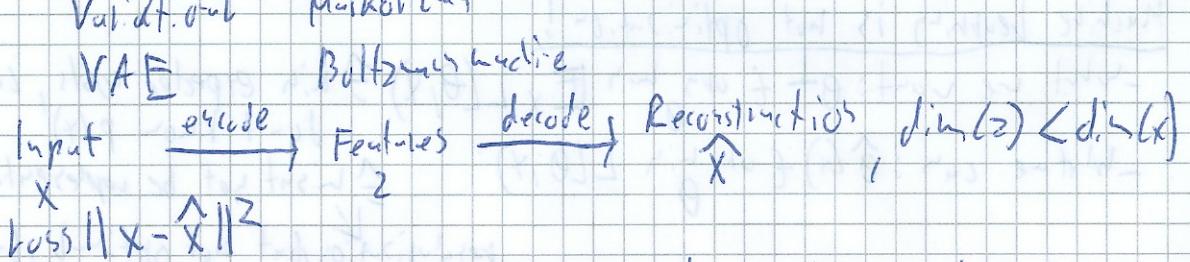
Markov chain
Generative stochastic
Networks (GSN)

Direct
Generative Adversarial
Network (GAN)

tractable density
- Fully visible belief nets
Pixel CNN

approximate
density

Autoencoder



For inference: throw away decoder and fixture or best funds

Denoising autoencoders: "blunk out" certain pixels in X

Want to learn invariants which are robust to translations, rotations, lighting

↳ augment data set
↳ regularization terms

↳ pre-process e.g. with SIFT

↳ local strategies e.g. CNNs

VAEs can generate data when choosing from representation

+ principled approach to generative models

+ allows inference of $q(z|x)$, can be useful for other tasks

- maximizes lower bound, not as good as more complex models like Pixel CNN

- submodels are blurrish and more complicated than GANs

GANs state-of-the-art, but hard to train, (can't use explicit density function) upsample network
 Generator network tries to fool discriminator by generating fake $g(z)$ image (NN)
 Discriminator network tries to distinguish between real and fake images, often (NN)
 $\min_{G} \max_{D} V(D, G) = \min_{G} \max_{D} \left\{ \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(g(z)))] \right\}$

for number of training iterations do

for k steps do

sample m noise samples (z^1, \dots, z^m) from noise prior $p_g(z)$

sample m examples $\{x^1, \dots, x^m\}$ from data distribution $p_{\text{data}}(x)$

update discriminator by ascending with $\nabla_D \frac{1}{m} \sum_{i=1}^m [\log D(x^i) + \log (1 - D(g(z^i)))]$
 end for

sample m noise samples from noise prior $p_g(z)$

update generator by descending with $\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m \log (1 - D(g(z^i)))$

end for

↳ in practice ascent on $\max_G \mathbb{E}_{z \sim p_z(z)} [\log (D(g(z)))]$

Machine learning is not optimization!

- what we want: $\theta^* = \arg \min_{\theta} \mathbb{E}_x L(\theta, x) \triangleq$ min expected costs, can't do because we don't know $p(x)$

- what we can: $\hat{\theta}^{(t)}$ fully with $L(\theta, x) \triangleq$ might not be representative, min empirical costs
 resolution, but not optimal, but typical

sensitivity analysis training $L(x') = \bar{x} + \Delta$, test $L(x'') = \bar{x} - \Delta$

$$\begin{aligned} \text{expected costs: } \mathbb{E}_{\theta|x'} [\log p(\theta|x')] &= \int p(\theta|\bar{x}+\Delta) \log p(\theta|\bar{x}-\Delta) d\theta \\ &= \int (p(\theta|\bar{x}) + \nabla_{\theta} p(\theta|\bar{x}+\Delta)|_{\theta=0} \cdot \Delta) \cdot (\log p(\theta|\bar{x}) - \frac{\nabla_{\theta} p(\theta|\bar{x}-\Delta)}{p(\theta|\bar{x}-\Delta)}|_{\theta=0} \cdot \Delta) d\theta \\ &= \mathbb{E}_{\theta|\bar{x}} [\log p(\theta|\bar{x})] - \Delta \int \nabla_{\theta} p(\theta|\bar{x}) (\log p(\theta|\bar{x}) - 1) d\theta + O(\Delta^2) \end{aligned}$$

Sample error is linear in Δ

$$\text{Posterior expectation: } \mathbb{E}_{\theta|x'} [p(\theta|x')] = \int p(\theta|\bar{x})^2 d\theta + O(\Delta^2)$$

↳ robustness

Lecture 12

k-means problem: find C and γ that minimize

$$R^{km}(c_1)) = \sum_{x \in X} \|x - M_{(c_1)}\|^2 \quad \text{cluster center}$$

k-means algorithm:

$\ln t \mu_t = x_i + \text{some random data points}$

for all x , $c(x) = \min_i \|x - \mu_i\|^2$ Expectation step.

for all $M / M_d = \frac{1}{n_d} \sum_{\{x\} = d} x$, $n_d = \# \text{att points in cluster}$ Maximization

until no changes in $\text{d}(x)$ are detected in practice.

EM for Mixture Gaussians

latent variable Soft assignment to component for every x and C
probabilistic of x in C

repeat

$$E: \gamma_{x_l} = \frac{PL(x_l, \theta^{(j)})}{PL(x|\theta^{(j)})}$$

$$M: M_L^{(j+1)} = \frac{\sum x_i \gamma_{xL}^{(j+1)}}{\sum \gamma_{xL}}, S_L^{2, j+1} = \frac{\sum \gamma_{xL} (x - M_L)^2}{\sum \gamma_{xL}} / P(c | \theta^{(j)})$$

until no changes in E-step

$$\text{solves } L(X | \pi_1, \dots, \pi_k, \theta_1, \dots, \theta_k) = \sum_x \log \sum_c \pi_c p(x|\theta_c)$$

had to track it log

amount of clusters is a hyperparameter

Beta distribution: $P(x|a,b) = \frac{1}{B(a,b)} \cdot x^{a-1} (1-x)^{b-1}, x \in [0,1], a,b > 0$

$$B(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}, \text{ where } \Gamma(u) = \int_0^{\infty} e^{-x} x^{u-1} dx$$

Dirichlet($x|\alpha$) = $\frac{1}{B(\alpha)} \prod_{k=1}^n x_k^{\alpha_k - 1}, x_i \in [0,1], \alpha_i > 0$

$$B(\alpha) = \frac{\prod_{k=1}^n \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^n \alpha_k)}$$

β dir is Dirichlet with only two variables

make "jelly"-mixing (lasts) $\hat{=}$ nonparametric Bayesian methods

Probabilities still have to sum up to 1

Use Dirichlet Process (α, H) — basic measure

st \hookrightarrow stick breaking construction

$$\beta_1 \sim \text{Beta}(1, \alpha), \rho_1 = \beta_1$$

$$\beta_2 \sim \text{Beta}(1, \alpha) \quad , \quad \rho_2 = \beta_2(1-\beta_1)$$

$$\beta_3 \sim \text{Beta}(1, \alpha) \quad , \quad \rho_3 = \beta_3(1-\beta_1)(1-\beta_2)$$

$$\beta_k \sim \text{Beta}(1, \alpha) \quad , \quad \rho_k = \beta_k(1-\sum_{i=1}^{k-1} \rho_i)$$

$G(\theta) = \sum_{k=1}^{\infty} \rho_k S_{\theta_k}(\theta) \sim DP(\alpha, H)$ \hookrightarrow GEM distribution

Dirac H

Chinese Restaurant Process

$$P(\text{customer } n+1 \text{ joins table } T | \mathcal{P}) = \begin{cases} \frac{|T|}{\alpha+n} & \text{if } T \in \mathcal{P} \\ \frac{\alpha}{\alpha+n} & \text{new table} \sim O(\text{Arry}(\alpha)) \end{cases}$$

\hookrightarrow process is older and labeling independent $\frac{\alpha}{\alpha+n}$ new table

\hookrightarrow exclusive!!!

Exchangability: (X_1, X_2, \dots) is a sequence of random variables. It is exchangeable if also under the permutation T (X_1, X_2, \dots) have the same distribution as $(X_{T(1)}, X_{T(2)}, \dots)$

De Finetti's theorem: $p(X_1, \dots, X_n) = \int \left(\prod_{i=1}^n p(x_i | G_i) \right) dP(G) \stackrel{\text{An i.i.d. exchangeable sequence}}{\leq} \stackrel{\text{can be represented by conditionally}}{\leq} \stackrel{\text{independent random variables}}{\leq}$

Pólya Urn: Pull a ball, put flat ball + one of same colour back, with

Hoppe Urn: Pull a ball, if coloured \rightarrow next colour, if black \rightarrow new colour, put back black and new colour

\hookrightarrow same as Chinese restaurant process

Fit DP Mixture Model with Gibbs Sampling

$$p(z_i=k | z_{-i}, x, \theta, \mu) \propto p(z_i=k | z_{-i}, \theta) p(x_i | z_i=k, z_{-i}, \mu)$$

$\hat{z}_i = \text{argmax}_k p(z_i=k | z_{-i}, \theta)$
(cluster assignment)

$\begin{cases} p(x_i | z_i=k, \mu) & \text{if } k \text{ exists} \\ p(x_i | \mu) & \text{if } k \text{ does not exist} \end{cases}$

Model selection by minimizing validation costs with posterior

$$p(\theta | X) = \exp(-\beta w(\theta, X)) \text{ with } w(\theta, X) = R(\theta, X) - F(X)$$

$$\text{Validation Loss: } \mathbb{E}_{\theta|x'} [-\log p(\theta|x')] = \beta \mathbb{E}_{\theta|x'} [w(\theta, x')] = \beta (\mathbb{E}_{\theta|x'} R(\theta, x') - F(x'))$$

posterior select. fn: with $\mathbb{E}_{\theta|x'} [-\log p(\theta|x')] \geq p(\cdot|.) - \log (\mathbb{E}_{\theta|x'} p(\theta|x'))$

$$= - \max_{\theta \in \Theta} \log \int p(\theta|x') p(\theta|x'') d\theta$$

probability kernel \rightarrow more similar \rightarrow better

Lecture 13

Generalization error: $R(\hat{c}) = P(\hat{c}(x) \neq c(x))$ $\rightarrow E[\hat{R}_n(c)] = R(c)$

Empirical error: $\hat{R}_n(c) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\hat{c}(x_i) \neq c(x_i)}$

Instance space $X \triangleq$ data in learners world

Concept (\hat{c}) \triangleq subset of X \subseteq classes

Concept class \triangleq a set of concepts we want to learn

Hypothesis class \triangleq another set of concepts that we use to learn the target concept

Learning algorithm \triangleq input: labeled samples $Z = \{(x_1, y_1), \dots, (x_n, y_n)\}$ with $y_i = c(x_i)$
output: $\hat{c} \in H$ -hypothesis class

Learnable: $\forall A$ can learn C if there is a poly-function such that

- for any distribution D on X and
- for any $0 < \epsilon < \frac{1}{2}$ and $0 < \delta < \frac{1}{2}$

A outputs \hat{c} such that $P_{Z \sim D^n}[R(\hat{c}) \leq \epsilon] \geq 1 - \delta$ / chance of being too wrong

with a sample Z of training set / how wrong can be
with a sample Z of size $\tilde{\gamma} \text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta}, \text{size}_C)$

PAC Learnable: concept class G is PAC learnable from a hypothesis class H
if there is an algorithm that can learn any concept in G

Efficient PAC Learnable: A learns in time polynomial in $\frac{1}{\epsilon}$ and $\frac{1}{\delta}$

The universal concept class $G = \{0, 1\}^*$ is not learnable from itself!

General setting: D on $X \times \{0, 1\}$ allows for classes having different labels

Learnable, instead of $R(\hat{c}) \leq \epsilon$ try to learn $(R(\hat{c}) - \inf_{c \in G} R(c)) \leq \epsilon$

General PAC Learning: - for any D on $X \times \{0, 1\}$ and

- for any $0 < \epsilon < \frac{1}{2}$ and $0 < \delta < \frac{1}{2}$

A receives sample Z of size $n \tilde{\gamma} \text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta}, \text{dim}(X))$

and can output \hat{c} such that

$P_{Z \sim D^n}[(R(\hat{c}) - \inf_{c \in G} R(c)) \leq \epsilon] \geq 1 - \delta$

Shattering: a set of instances A can be shattered by a concept class \mathcal{C}
if for every subset $S \subseteq A$, there is a concept $h \in \mathcal{C}$ such that $h(S) = S$

e.g. set of two points in \mathbb{R}^2 and $\mathcal{C} \subseteq$ rectangles
but 3 numbers in \mathbb{R} and $\mathcal{C} \subseteq$ intervals

VL-dimensions: dimensions needed to shatter set \rightarrow complexity measure for \mathcal{C}

Weak learner: PAC learnable only for "hard" errors ϵ
Strong learner: PAC learnable for all errors ϵ

Error probability for realizable finite hypothesis classes $H = \mathcal{C}$

Assume $\min_{\mathcal{C}} R(c) = 0$, $\hat{c} \in \arg\min_{\mathcal{C}} \hat{R}_n(c)$

$$\begin{aligned} P(R(\hat{c}) > \epsilon) &= P \left(\max_{\mathcal{C}, \hat{R}_n(c)=0} R(c) > \epsilon \right) \quad \text{max deficiency in sum} \\ &\leq \sum_{\mathcal{C}} P(R(c) > \epsilon \wedge \hat{R}_n(c) = 0) \\ &\leq |\mathcal{C}| (1 - R(c))^n \leq |\mathcal{C}| e^{-n\epsilon} \leq \delta \\ &\leq n \cdot \frac{1}{\epsilon} (\log |\mathcal{C}| + \log \frac{1}{\delta}) \end{aligned}$$

Vapnik-Chervonenkis inequality

$\hat{c} \in \arg\min_{\mathcal{C}} \hat{R}_n(c) = \arg\min_{\mathcal{C}} |\{x_i : \hat{c}(x_i) \neq y_i\}|$ = minimizer of empirical risk

$c^* \in \arg\min_{\mathcal{C}} R(c) \approx$ minimizer of expected risk

$$R(\hat{c}) - \inf_{\mathcal{C}} R(c) = R(\hat{c}) - \hat{R}_n(\hat{c}) + \hat{R}_n(\hat{c}) - R(c^*) \leq 2 \sup_{\mathcal{C}} |R(c) - \hat{R}_n(c)|$$