

Gradient Descent

feature number
↓

Wish to fit the model $h(x) = \sum_{i=0}^n \theta_i x_i$ (let $x_0 = 1$)

- Define cost function $J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$ [least-square cost function] [Reason to choose: $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$]
 \downarrow
 $i^{\text{th}} \text{ evaluation}$

(IID) $\epsilon^{(i)} \sim N(0, \sigma^2)$, $y^{(i)} | x^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$

minimize $J \Rightarrow \text{MLE to find } \hat{\theta}$

- Gradient Descent. $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$

↓
(learning rate)

- where $\frac{\partial}{\partial \theta_j} \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 = \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \Rightarrow \theta_j := \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$, for every j

$$\Rightarrow \theta := \theta + \alpha \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x^{(i)} \quad [\text{batch gradient-descent}]$$

↓
number of data

- choice of learning rate α : too small — takes too long for convergence
too big — takes too large step and run past minimum
(strong signal; $J(\theta)$ increasing instead of decreasing)

- Stochastic Gradient Descent (for large data-set)

for $i=1$ to n :

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} \quad [\text{for one sample do improvement}]$$

(for every i)

- Note: this algorithm will not self-converge like batch-gradient descent (i.e. not using parameters from the global optimum)
- when using stochastic gradient descent, decrease the learning rate

Normal Equation

(for gradient descent to reach global optimum in "one" step)

Solve normal equation $\nabla_{\theta} J(\theta) = \vec{0}$. (Eq. 2.1)

Define Design Matrix X (of samples) = $\begin{pmatrix} -(x^{(1)})^T \\ \vdots \\ -(x^{(m)})^T \end{pmatrix}$ → train sample vector transposed.

$$\text{then } X\theta = \begin{pmatrix} x_0^T \\ x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} = \begin{pmatrix} (h_\theta(x^{(1)}))^T \\ (h_\theta(x^{(2)}))^T \\ \vdots \\ (h_\theta(x^{(m)}))^T \end{pmatrix} \Rightarrow X\theta - \vec{y} = \begin{pmatrix} (h_\theta(x^{(1)}) - y^{(1)})^T \\ (h_\theta(x^{(2)}) - y^{(2)})^T \\ \vdots \\ (h_\theta(x^{(m)}) - y^{(m)})^T \end{pmatrix}$$

$$\text{then } \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 = J(\theta)$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \nabla_{\theta} (\theta^T X^T X \theta - (\vec{y})^T (X\theta) + \vec{y}^T \vec{y}) \quad (\text{Eq. 2.2})$$

Now Notice: $J: \mathbb{R}^{mn} \rightarrow \mathbb{R}$, generalize: for $f: \mathbb{R}^{mn} \rightarrow \mathbb{R}$, define $\nabla_A f(A) = \begin{pmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{n1}} & \cdots & \frac{\partial f}{\partial A_{nn}} \end{pmatrix}$

$$\text{where } A = \begin{pmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{n1} & \cdots & A_{nn} \end{pmatrix}$$

Claim 1. $f(A) = \text{tr}AB$, where $A, B \in \mathbb{R}^{n \times n}$, $f(A) = \sum_{i=1}^n \sum_{j=1}^n A_{ij}B_{ji} \Rightarrow \nabla_A f(A) = B^T$

Claim 2. $\text{tr}ABC = \text{tr}CAB$.

Claim 3. $\nabla_A \text{tr}AATC = CA + C^TA$.

$(\text{Eq 2.2}) \Rightarrow \frac{1}{2} \nabla_\theta (\theta^T X^T \theta - (\vec{y})^T (X\theta) - (X\theta)^T (\vec{y}))$

$= \frac{1}{2} \nabla_\theta (\theta^T X^T \theta - (\vec{y})^T (X\theta) - (\vec{y})^T (X\theta))$

$= \frac{1}{2} \nabla_\theta (\theta^T X^T \theta - 2(\vec{X}\vec{y})^T \theta)$

$= \frac{1}{2} \nabla_\theta (\text{tr}(\vec{X}\vec{X})\theta\theta^T - 2\text{tr}(\theta(\vec{X}\vec{y})^T))$

$\nabla_\theta \text{tr}(\theta\theta^T) \leftarrow \text{Claim 2+3} \rightarrow \nabla_\theta \text{tr}(\theta(\vec{X}\vec{y})^T)$

$= \frac{1}{2} [X^T X \theta + (\vec{X}\vec{X})^T \theta - 2\vec{X}\vec{y}]$

$= X^T X \theta - \vec{X}\vec{y} \xrightarrow{\text{set } \vec{\theta}} \vec{\theta} = (\vec{X}\vec{X})^{-1}(\vec{X}\vec{y})$ (if $\vec{X}\vec{X}$ is invertible, then calc $\vec{\theta}$ directly)

Locally Weighted Regression

Adjustments as compared to linear regression:

Fit θ to minimize $\sum_{i=1}^m w^{(i)}(y^{(i)} - \theta^T x^{(i)})^2$, $w^{(i)}$ is a weighting function defaulted by: $w^{(i)} = \exp\left(-\frac{(x^{(i)} - \bar{x})^2}{2\tau^2}\right)$

[intuitively, if $|x^{(i)} - \bar{x}|$ is small, $w^{(i)} \approx 1$; if $|x^{(i)} - \bar{x}|$ is big, $w^{(i)} \approx 0$]

- Sensitivity test for the selection of τ
- Use locally weighted regression when feature numbers are small (i.e. $\geq 3, 4$)
- Can choose τ according to the impression held by Gaussian Distribution (i.e. $\frac{1}{\sqrt{2\pi}\tau} e^{-\frac{(x-\mu)^2}{2\tau^2}}$)

Logistic Regression

want $h_\theta(x) \in [0, 1] \Rightarrow$ choose $g(z) = \frac{1}{1+e^{-z}}$ s.t. to a hypothesis $h_\theta(x)$, $h_\theta(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$

("sigmoid" / "logistic" function)

here $y \in \{0, 1\}$, for conventional reasons, choose (assume)

$$P(y=1 | x; \theta) = h_\theta(x) \Rightarrow P(y=0 | x; \theta) = h_\theta(x) [1 - h_\theta(x)]^{1-y}$$

$$P(y=0 | x; \theta) = 1 - h_\theta(x)$$

$$y \in \{0, 1\}$$

Assuming that n training samples were generated independently, the likelihood of the parameters:

$$L(\theta) = P(\bar{y} | x; \theta) = \prod_{i=1}^n (h_\theta(x^{(i)})^{y^{(i)}}) (1 - h_\theta(x^{(i)}))^{1-y^{(i)}}$$

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1 - h_\theta(x^{(i)})) \quad [\text{Choose } \theta \text{ to maximize } \ell(\theta)]$$

- gradient ascent: $\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} \ell(\theta)$ (\leftarrow when deriving this part remember $h_\theta(x)$ is NON-LINEAR)

$$\begin{aligned} \text{where } \frac{\partial}{\partial \theta_j} \ell(\theta) &= \left(y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left(y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) g(\theta^T x) (1-g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1-g(\theta^T x)) - (1-y)g(\theta^T x)) x_j = (y - h_\theta(x)) x_j \end{aligned}$$

- stochastic gradient ascent updating: $\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$ [$L(\theta)$ guaranteed concave \Rightarrow only global optimum]

? (Q: How to set the boundary since it won't self-converge) \downarrow
 + Consider Newton's method for solving $\ell'(\theta) = 0$ (i.e. $\theta^{(t+1)} = \theta^{(t)} - \frac{\ell'(\theta^{(t)})}{\ell''(\theta^{(t)})}$ if $\theta^{(t)} \in \mathbb{R}^n$) \downarrow decrease learning rate until θ doesn't change
 (quadratic convergence)

- When θ is a vector (say $\theta \in \mathbb{R}^n$) $\theta^{(t+1)} := \theta^{(t)} + H^{-1} \nabla \ell$

where: $\nabla \ell$ is a vector of derivative $\ell \in \mathbb{R}^n$

$H \in \mathbb{R}^{(n+1) \times (n+1)}$ and $H_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j}$,
 ↓ Hessian \downarrow every iteration of Newton can be expensive.

but Newton's Method quickly converges,

when calculating H^{-1} .

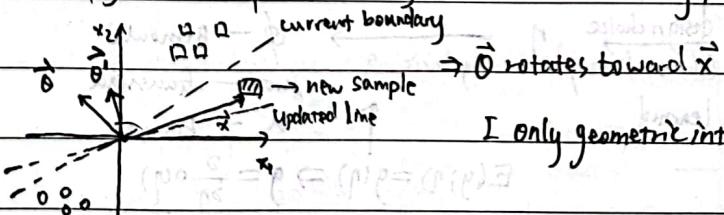
Perceptron Learning

Similar to logistic regression, we now choose $g(z) = \begin{cases} 1, z \geq 0 \\ 0, z < 0 \end{cases}$

and $h_\theta(x) = g(\theta^T x)$, with an updating rule $\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$ [this is different from logistic]

Notice: $y^{(i)} - h_\theta(x^{(i)}) = 0 \Rightarrow$ right prediction

$|y^{(i)} - h_\theta(x^{(i)})| > 1 \Rightarrow$ wrong prediction, a wrong "prediction" contributes to update;



I only geometric interpretation, no probabilistic interpretation).

Generalized Linear Models

- Exponential family $P(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$ [given canonical parameters]

where • often the case that $T(y) = y$

• $a(\eta)$ plays a role in normalization (i.e. $\int b(y) \frac{\exp(\eta^T T(y))}{\exp(a(\eta))} dy = 1$)

• $\eta, T(y)$ must match in dimension

given a fixed choice of $T(y)$, $a(\eta)$ and $b(y)$ \Rightarrow vary η to get different distributions in this exponential family.

- can also include dispersion parameter s.t. $p(y; \eta, \tau) = b(a, \tau) \exp(\eta^T T(y) - \frac{a(\eta)}{c(\tau)})$

in Gaussian, $c(\tau) = \sigma^2$

(negative log likelihood)

- Properties \oplus MLE wrt to $\eta \Rightarrow$ concave; NLL wrt $\eta \Rightarrow$ convex.

$$\textcircled{2} E(y; \eta) = \frac{\partial}{\partial \eta} a(\eta)$$

$$\textcircled{3} \text{Var}(y; \eta) = \frac{\partial^2}{\partial \eta^2} a(\eta)$$

[if η is a vector, then operator would be a Hessian]

[-- Choices: Real (numbers) - Gaussian

Binary - Bernoulli

Count (positive int) - Poisson

ii) $y|x; \theta \sim \text{Exponential Family}(\eta)$ [given x and θ , $y \sim \text{EF}(\eta)$]

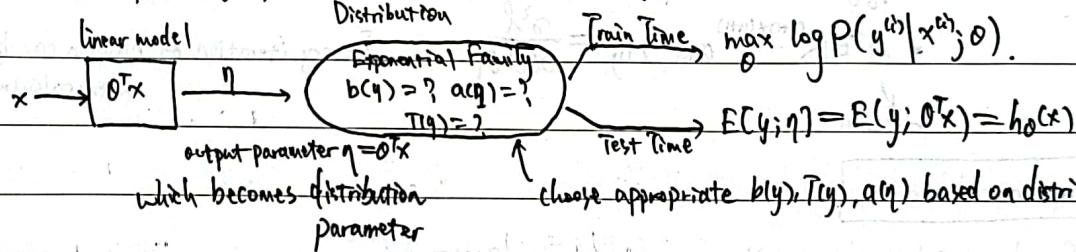
R^+ (positive real) - Gamma, Exponential

12) Given x , $\eta = \theta^T x$

[if η is vector then $\eta_i = \theta_i^T x$]

13) At Test time, output $E(y|x; \theta)$

$[h_\theta(x) = E(y|x; \theta)]$



Incidental Benefits:

1) learning update rule: $\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$ [Only difference is the distribution $h_\theta(x)$].

2) Three (kinds) of Parameterizations

| Model Parameter | Natural Parameter | Canonical Parameter |
|--------------------------------|---|---------------------|
| θ | η <small>design choice (choose linear)</small> | ϕ - Bernoulli |
| \uparrow | η <small>algebraic</small> | μ - Gaussian |
| the only thing that is learned | λ - Poisson | |

$$E(y; \eta) = g(\eta) \Rightarrow g = \frac{\partial}{\partial \eta} a(\eta)$$

① Choose distribution for y

13) Ex. in logistic regression, ** The choice distribution of y is Bernoulli **

② Calculate Expectation of distribution ①

$$\Rightarrow h_\theta(x) = g(\theta^T x) = E(y|\eta) = \phi, \text{ here in Bernoulli, } P(y|\phi) = \exp((\log \frac{\phi}{1-\phi})y + \log(1-\phi))$$

③ Write in terms of Exp Family to calculate $g = \frac{\partial}{\partial \eta} a(\eta)$

$$\Rightarrow a(\eta) = \log(1-\phi) \Rightarrow \eta = \frac{\partial}{\partial \eta} a(\eta) = \frac{1}{1+\phi} = \frac{1}{1+e^{-\eta}}$$

$$= \log(1+e^{-\eta})$$

$$\Rightarrow \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}} \quad \textcircled{4} \text{ Equation } E(y|x; \theta) = g(\eta) \text{ to solve for canonical parameters}$$

SoftMax Regression

$\phi_{ik} = \frac{1}{\sum_i} \phi_i$ since $(k-1)$ independent variables

Classification Problem with $y \in \{1, 2, \dots, k\}$, use $\phi_1, \phi_2, \dots, \phi_k$ to denote possibility of belonging to class i .

define $T(y) \in \mathbb{R}^{k+1}$ as one-hot vector, where:

$$T(1) = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, T(2) = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \dots, T(k) = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix} \text{ and } T(k) = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\text{obtain } E(T(y))_i = P(y=i) = \phi_i$$

$[T(y)_i]$ denotes i -th element in $T(y)$

- then $P(y; \theta) = \phi_1 \phi_2 \dots \phi_k$

ExpF. Interpretation
 $= \exp \left(T(y)_1 \log \phi_1 + \dots + T(y)_k \log \phi_k + \left(1 - \sum_{i=1}^k T(y)_i \right) \log \phi_k \right)$
 $= \exp \left(\begin{pmatrix} \log \phi_1 / \phi_k \\ \vdots \\ \log \phi_k / \phi_k \\ \vdots \\ T(y)_k \end{pmatrix}^T \begin{pmatrix} T(y)_1 \\ T(y)_2 \\ \vdots \\ T(y)_k \end{pmatrix} + \log \left(1 - \sum_{i=1}^k \phi_i \right) \right)$ [belongs to Exp Family]

the link function is given by $\eta_i = \log \phi_i / \phi_k \Rightarrow \phi_i = \frac{1}{\sum_j e^{\eta_j}}$ [softmax function]

given η 's and x 's are linearly linked, $\eta_i = \theta_i^T x$ ($i = 1, 2, \dots, k-1$) and $\theta_i \in \mathbb{R}^{d+1}$ then.

$$p(y|x; \theta) = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}}$$

output is also a distribution

$$h_\theta(x) = E[T(y)|x; \theta] = (\phi_1, \phi_2, \dots, \phi_{k-1})^T = \left(\frac{e^{\theta_1^T x}}{\sum_{j=1}^k e^{\theta_j^T x}}, \dots, \frac{e^{\theta_{k-1}^T x}}{\sum_{j=1}^k e^{\theta_j^T x}} \right)$$

[output probability $p(y=i|x; \theta)$]

log-likelihood given n samples $\{x^{(i)}, y^{(i)}\}$ ($i = 1, 2, \dots, n$)

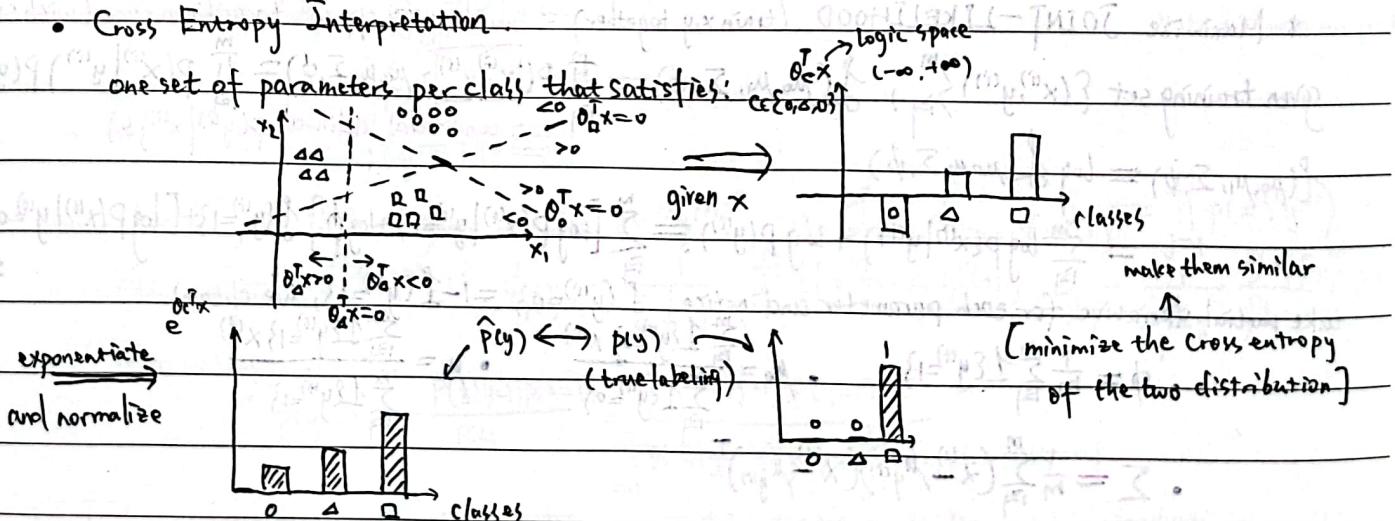
$$\ell(\theta) = \sum_{i=1}^n \log p(y^{(i)}|x^{(i)}; \theta) = \sum_{i=1}^n \log \prod_{j=1}^k \left(\frac{e^{\theta_j^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)$$

[$1 \cdot 3 = 1$ if \cdot else 0]

maximize $\ell(\theta) \Leftarrow$ gradient ascent or Newton's method.

Cross Entropy Interpretation

one set of parameters per class that satisfies:



$\hat{p}(y)$ for given $x \rightarrow$ a distribution

(likelihood of given x belonging to which class)

$$\text{Cross Entropy } (\hat{p}(y), p(y)) = -\sum_{x \in \text{label}} p(y) \log \hat{p}(y) = -\sum_{\substack{y \in \{0, 1, 0\} \\ \text{p(y) being ideal/true}}} \left(\begin{array}{l} p(y) = 1 \text{ for } 0 \\ p(y) = 0 \text{ for } 1 \\ p(y) = 0 \text{ for } 0 \end{array} \right) \log \hat{p}(y)$$

$$= -\log \hat{p}(y) = -\log \frac{e^{\alpha x}}{\sum_{y \in \{0, 1, 0\}} e^{\alpha x}}$$

label

\Rightarrow do gradient descent to minimize cross entropy.

Generative Learning Algorithm

Idea: given y , what do features look like? (namely learn $p(x|y)$ as opposed to discriminative learning which learns $p(y|x)$)
say that $y \in \{0, 1\}$, according to Bayes: $P(y=1|x) = \frac{P(x|y=1)p(y=1)}{P(x|y=1)p(y=1) + P(x|y=0)p(y=0)}$

learned determined by distribution of y

Gaussian Discriminant Analysis (GDA)

* Assume: given y , features x is distributed according to multivariate normal distribution.

* $\mathbf{z} \sim N(\mu, \Sigma)$, where $\mathbf{z} \in \mathbb{R}^n$, $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ being the covariance matrix.

where $\mathbf{E}(\mathbf{z}) = \mu$

$\bullet \text{Cov}(\mathbf{z}) = \mathbf{E}((\mathbf{z}-\mu)(\mathbf{z}-\mu)^T) = \mathbf{E}[\mathbf{zz}^T] - (\mathbf{E}(\mathbf{z}))(\mathbf{E}(\mathbf{z}))^T = \Sigma$

$\bullet P(\mathbf{z}) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2}(\mathbf{z}-\mu)^T \Sigma^{-1}(\mathbf{z}-\mu))$.

Based on the assumed distribution, GDA model:

$$p(x|y=1) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)) \quad \begin{matrix} \text{parameters: } \mu_0, \mu_1, \Sigma, \phi \\ \text{different means for each classes } \in \mathbb{R}^n \end{matrix}$$

$$p(x|y=0) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1}(x-\mu_0)) \quad \begin{matrix} \text{same cov for each } \in \mathbb{R}^n \\ \text{can make the covariance matrix different} \end{matrix}$$

with y considering Bernoulli distributed: $p(y) = \phi^y (1-\phi)^{1-y}$ for each class, but this would give a non-linear

* Maximize JOINT-LIKELIHOOD. (train x, y together) boundary as opposed to a linear one (with same Σ)

$$\text{given training set } \{(x^{(i)}, y^{(i)})\}_{i=1}^m \quad \ell(\mu_0, \mu_1, \Sigma, \phi) = \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \mu_0, \mu_1, \Sigma, \phi) = \prod_{i=1}^m p(x^{(i)}|y^{(i)}) p(y^{(i)})$$

$$\ell(\mu_0, \mu_1, \Sigma, \phi) = \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \mu_0, \mu_1, \Sigma, \phi)$$

$$= \sum_{i=1}^m [\log p(x^{(i)}|y^{(i)}) + \log p(y^{(i)})] \quad \begin{matrix} \text{not conditional likelihood } p(y^{(i)}|x^{(i)}; \theta) \\ \text{but joint likelihood } p(x^{(i)}, y^{(i)}; \theta) \end{matrix}$$

take partial derivative for each parameter and notice $\int_{\mathbb{R}^n} 1_{\{y^{(i)}=0\}} = 1 - \int_{\mathbb{R}^n} 1_{\{y^{(i)}=1\}}$, we obtain:

$$\bullet \phi = \frac{1}{m} \sum_{i=1}^m 1_{\{y^{(i)}=1\}} \quad \bullet \mu_0 = \frac{1}{m} \sum_{\substack{i=1 \\ 1_{\{y^{(i)}=0\}}}}^m x^{(i)} \quad \bullet \mu_1 = \frac{1}{m} \sum_{\substack{i=1 \\ 1_{\{y^{(i)}=1\}}}}^m x^{(i)}$$

$$\bullet \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

when making predictions, model output $\arg\max_y p(y|x) = \arg\max_y p(y)p(x|y)$.

Linkage with Logistic Regression

If plotting $p(y=1|x; \phi, \mu_0, \mu_1, \Sigma)$ as a function of x , we'll find the form

$$p(y=1|x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1+e^{-\phi x}} \text{ for some appropriate } \phi.$$

GDA assumptions (generative) Logistic Regression assumptions (discriminative)

$$\begin{aligned} x|y=0 &\sim N(\mu_0, \Sigma) \\ x|y=1 &\sim N(\mu_1, \Sigma) \\ y &\sim \text{Bernoulli}(\phi) \end{aligned} \quad \xrightarrow{\text{indicates}} \quad p(y=1|x) = \frac{1}{1+e^{-\phi x}} \quad x|y \sim \text{Exp Family}$$

and for $y=k$ and $y=j$, if $k \neq j$ then
 $x|y$ only vary in natural parameters

• GDA - stronger set of assumptions; LR - weaker set of assumptions $\Rightarrow p(y|x)$ is logistic

if assumptions are roughly correct,
then model produced ~~won't do better~~ \downarrow weaker assumptions make the model

more robust, a safe choice if don't know
(even with a small dataset) $\quad \quad \quad$ the distribution of $x|y$.

Naive Bayes

* Turn occurrence of word into feature vectors $\vec{x} = \{0, 1\}^n$, where n is the total number of feature words

and $x_i = 1$ {Word i appears in current piece of text}

* Want to model $p(x|y)$ and $p(y)$, direct modeling using $2^n - 1$ parameters won't work

* Assume x_i 's are conditionally independent given y , then

$$\begin{aligned} p(x_1, x_2, \dots, x_n | y) &= p(x_1|y)p(x_2|y)p(x_3|y) \dots p(x_n|y) \quad \text{extremely strong} \\ &= p(x_1|y)p(x_2|y) \dots p(x_n|y) \quad [\text{Naive Bayes Assumption (not exactly true)}] \end{aligned}$$

Parametrized by

- $\phi_{j|y=1} = P(x_j=1 | y=1)$ \downarrow but works well on problems
- $\phi_{j|y=0} = P(x_j=0 | y=0)$
- $\phi_y = P(y=1)$

Then the joint likelihood given samples $\{x^{(1)}, y^{(1)}\}, \dots, \{x^{(m)}, y^{(m)}\}$ is $L(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{j=1}^m p(x_j^{(i)}, y^{(i)})$
MLE: $\phi_y = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y^{(i)}=1\}$, $\phi_{j|y=1} = \frac{1}{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=1\}} \sum_{i=1}^m \mathbb{1}\{x_j^{(i)}=1 \wedge y^{(i)}=1\}$, $\phi_{j|y=0} = \frac{1}{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=0\}} \sum_{i=1}^m \mathbb{1}\{x_j^{(i)}=0 \wedge y^{(i)}=0\}$

when making predictions,

$$P(y=1|x) = \frac{P(x|y=1)p(y=1)}{P(x)} = \frac{\prod_{j=1}^n P(x_j|y=1)p(y=1)}{\prod_{j=1}^n P(x_j|y=1)p(y=1) + \prod_{j=1}^n P(x_j|y=0)p(y=0)}$$

[Note]: when x_j can take values of $\{1, 2, \dots, k\}$, can model $p(x_j|y)$ as multinomial instead of Bernoulli.

Moreover, can discretize continuous values to use Naive Bayes when GDA didn't turn out well.

e.g. x_i \downarrow < 400 \downarrow $400-800$ \downarrow $800-1200$ \downarrow > 1200 \downarrow 1 \downarrow 2 \downarrow 3 \downarrow 4 \downarrow [conventionally, discretize into 10 values]

↓ Extension: Say maybe a word didn't exist in the training set, then with MLE some of ϕ_j might end up as zero.

Apply Laplace Smoothing: $\hat{\phi}_j = \frac{1}{k+n} \left(\sum_{i=1}^n \mathbb{1}\{x_i^{(i)} = j\} + 1 \right)$, where $\sum_{j=1}^k \hat{\phi}_j = 1$ still holds. [given based on Bayesian - Estimate].

Then the estimate for parameters becomes:

$$\hat{\phi}_{j|y=1} = \frac{1 + \sum_{i=1}^m \mathbb{1}\{x_i^{(i)} = j, y_i^{(i)} = 1\}}{2 + \sum_{i=1}^m \mathbb{1}\{y_i^{(i)} = 1\}}, \quad \hat{\phi}_{j|y=0} = \frac{1 + \sum_{i=1}^m \mathbb{1}\{x_i^{(i)} = j, y_i^{(i)} = 0\}}{2 + \sum_{i=1}^m \mathbb{1}\{y_i^{(i)} = 0\}}, \quad \hat{\phi}_y = \frac{\sum_{i=1}^m \mathbb{1}\{y_i^{(i)} = 1\}}{m}$$

↑ whether or not applying LS
to ϕ_y doesn't matter too much.

Support Vector Machine

• Functional Margin and Geometric Margin

consider binary classification problem with labels $y \in \{-1, 1\}$, and rewrite classifier as

$$h_{w,b}(x) = g(w^T x + b), \text{ where } g(z) = \begin{cases} +1, & \text{if } z \geq 0 \\ -1, & \text{otherwise.} \end{cases} \quad \underbrace{[b_0, b_1, b_2, \dots, b_d]}_b \quad \underbrace{w^T x + b}_w$$

* Functional Margin. (drop $x_0 = 1$ convention)

- Given a training example $(x^{(i)}, y^{(i)})$, define functional margin of w, b as

$$\hat{y}^{(i)} = y^{(i)}(w^T x^{(i)} + b), \quad \begin{cases} \text{if } y^{(i)} = 1, \text{ want } w^T x^{(i)} + b > 0 \\ \text{if } y^{(i)} = -1, \text{ want } w^T x^{(i)} + b < 0 \end{cases} \quad \text{if } \hat{y}^{(i)} > 0, \text{ then right prediction.}$$

[Notice] if scaling $(w, b) \rightarrow (kw, kb)$ ($k > 1$) ↓ confidence correctness.

we can make $\hat{y}^{(i)}$ arbitrarily large without effecting the sign of $\hat{y}^{(i)}$

⇒ impose normalization condition such as $\|w\|_2 = 1$ or $(w, b) \rightarrow \left(\frac{w}{\|w\|_2}, \frac{b}{\|w\|_2}\right)$.

- Given a training set, want the functional margin $\hat{Y} = \min_i \hat{y}^{(i)}$ to be big.

* Geometric Margin

$$\vec{OB} = \vec{OA} - \vec{x}^{(i)}, \quad \vec{w} = \vec{x}^{(i)} - \vec{y}^{(i)} \frac{w}{\|w\|_2} = \vec{x}^{(i)} - \vec{y}^{(i)} \frac{w}{\|w\|_2}$$

where $w^T(x_0) + b = 0 \Rightarrow \vec{y}^{(i)} = \frac{w^T x^{(i)} + b}{\|w\|_2}$

$$\text{more generally, define gm wrt } (x^{(i)}, y^{(i)}) \quad \hat{y}^{(i)} = y^{(i)} \left(\frac{w}{\|w\|_2} \right)^T x^{(i)} + \frac{b}{\|w\|_2} = \frac{\hat{y}^{(i)}}{\|w\|_2}$$

can impose arbitrary scale on $\hat{y}^{(i)}$ [invariant to scaling of parameters]

Define geometric margin of $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ to be

$$\gamma = \min_i \hat{y}^{(i)}$$

* Optimal Margin Classifier

Assume the training set is linearly separable, pose optimization problem:

$$\max_{w, b} \gamma$$

[non-convex optimization] (1)

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \gamma, \quad i = 1, 2, \dots, n.$$

$$\|w\|_2 = 1$$

Notice we can add arbitrary scaling constraint to w, b , and $(kw, kb) \rightarrow k\hat{r}$, we can assume $\hat{r} = 1$

$$\text{w.r.t training set, and farther notice } \max_{\substack{\|w\| \\ \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \hat{Y}}} \hat{Y} \iff \min_{\substack{w, b \\ \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, i=1,2,\dots,n}} \frac{1}{2} \|w\|^2 (\|w\|)$$

which is, switch \hat{Y} in Eq(1) to \hat{Y} , and choose $\|w\| = \frac{1}{2}$ to get Eq(2). [quadratic objective and linear constraints]

↓ Extension: Lagrange Duality

Consider primal optimization: $\min_w f(w)$

$$\text{s.t. } g_i(w) \leq 0, i=1,2,\dots,k$$

$$h_i(w) = 0, i=1,2,\dots,l$$

Define generalized Lagrangian $\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$ [α_i, β_i being lagrange multipliers]

* PRIMAL * $\mathcal{O}_P(w) = \max_{\alpha, \beta \geq 0} \mathcal{L}(w, \alpha, \beta) = \begin{cases} f(w), & \text{if } w \text{ satisfies primal constraints} \\ \infty, & \text{o.w.} \end{cases}$

then $p^* \triangleq \min_w \mathcal{O}_P(w) = \min_w \max_{\alpha, \beta \geq 0} \mathcal{L}(w, \alpha, \beta)$ is the same optimal problem as (3)

* DUAL * $d^* \triangleq \max_{\alpha, \beta \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) = \max_{\alpha, \beta \geq 0} \mathcal{O}_D(\alpha, \beta)$

We obtain that

$$d^* = \max_{\alpha, \beta \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta \geq 0} \mathcal{L}(w, \alpha, \beta) = p^*$$

(Under Conditions of: • f and g_i 's are convex, and h_i 's are affine (linear with intercept))

if f has a Hessian, then f is convex $\Leftrightarrow H_f$ is positive + semidefinite

• constraints g_i are strictly feasible $\Rightarrow \exists w \in \mathbb{R}^n$ s.t. $g_i(w) < 0$ for $i \in \{1, 2, \dots, k\}$.

Then there must exist w^*, α^*, β^* s.t. w^* is the solution to the primal problem; α^*, β^* are solutions

to the dual problem, and moreover $p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$.

[Note: ① w^*, α^*, β^* satisfy KKT

Moreover, w^*, α^*, β^* satisfy the KKT conditions;

② they are solutions for optimality

$$\cdot \frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0 \quad \text{③ if } p^* = d^* \text{ and } w^*, \alpha^*, \beta^* \text{ is solution}$$

⇒ they must satisfy KKT.

$$\cdot g_i(w^*) \leq 0, \alpha^* \geq 0$$

(dual complementarity) $\alpha_i^* g_i(w^*) = 0, i=1,2,\dots,k$.

[also: if some w^*, α^*, β^* satisfy KKT ⇒ solution to

primal & dual

↑ Optimal Margin Classifier

points closest to the separating hyperplane have

$$g_i(w) = -y^{(i)}(w^T x^{(i)} + b) + 1 = 0 \quad [\text{rather than } \leq 0]$$

Hence, only these points hold $\alpha_i^* > 0$

[support vectors]

$$-y^{(i)}(w^T x^{(i)} + b) + 1 = 0$$

max margin separating hyperplane

Construct Lagrangian $\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)}(w^T x^{(i)}) + b] - 1$ [solve the dual problem].

• Find $\min_{w,b} \mathcal{L}(w, b, \alpha)$ which is $\mathcal{O}_D(\alpha)$. $\Rightarrow \nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0 \Rightarrow w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$.

$$\Rightarrow \frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

[w lies in the span of training sets].

• Then $\mathcal{L}(w, b, \alpha)$ simplifies to $\mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (x^{(i)})^T x^{(j)} - b \sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle$

$$= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle$$

• Get the following optimization problem:

$$\max_{\alpha} N(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle$$

s.t. $\alpha_i \geq 0, i=1, 2, \dots, n$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0.$$

(satisfy conditions for $p^* = d^*$ and holds KKT conditions]

• if we can solve for optimal α in this case, then

$$w^* = \sum_{i=1}^n \alpha_i^* y^{(i)} x^{(i)}$$

• Consider the primal problem $\min_{w,b} \mathcal{L}(w, b, \alpha) = \min_{w,b} \mathcal{O}_D(w, b) = \min_{w,b} \left\{ \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1] \right\}$

[consider $g(w) = -y^{(i)}(w^T x^{(i)} + b) + 1$]

plug in solutions of dual: $w^* = \sum_{i=1}^n \alpha_i^* y^{(i)} x^{(i)}$ $b = w^T x_k - \bar{y}_k$ for $\alpha_k > 0$ \hookrightarrow (Q: mathematical instead?)

with $y^{(i)} = 1 \Rightarrow b \geq -(w^T x^{(i)}) + 1; y^{(i)} = -1 \Rightarrow b \leq -w^T x^{(i)} - 1$, geometrically maximize geometric margin of both side.

$$\Rightarrow b^* = -\frac{1}{2} \left[\max_{i:y^{(i)}=1} (w^* T x^{(i)}) + \min_{i:y^{(i)}=-1} (w^* T x^{(i)}) \right] = b^* = -\frac{1}{2} \left[\max_{i:y^{(i)}=1} (w^* T x^{(i)}) + \min_{i:y^{(i)}=-1} (w^* T x^{(i)}) \right].$$

Notice when making a new prediction, $w^T x + b = \left(\sum_{i=1}^n \alpha_i y^{(i)} \langle x^{(i)}, x \rangle \right) + b$ [calculate quantity of inner product between x_{test} and x_{train}].

Moreover, based on KKT conditions, α_i would be non-zero for support vectors.

\Rightarrow only need to find inner product between x and the support vectors. [of which there is often only a small number]

Kernel Methods

* Define kernel corresponding to feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^P$ as a function that maps $X \times X \rightarrow \mathbb{R}$ satisfying

$$K(x, z) \triangleq \langle \phi(x), \phi(z) \rangle \quad [\text{in this case, } X = \mathbb{R}^d]$$

Logic behind kernels: if algorithm can be expressed in terms of $\langle x^{(i)}, x^{(j)} \rangle$, but $x^{(i)} \in \mathbb{R}^d$ where d is (and need to be done repeatedly)

too large thus calculating $\langle x^{(i)}, x^{(j)} \rangle$ in $O(d)$ time can be costly, we can map $x^{(i)}$ to higher dimensional

feature $\phi(x^{(i)})$ where $K(x, z)$ is precomputed and saves time in calculation and then work implicitly in

the mapping space \mathbb{R}^P without explicitly referring to original $x^{(i)}$ in \mathbb{R}^d dimension.

Example] Least mean square with kernels.

In section of linear regression, we need to minimize $J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$, to do so we refer to

batch gradient descent update: $\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)})) x^{(i)}$ $\xrightarrow{\text{linear model}}$ $\theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) x^{(i)}$

Let $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^P$, now our goal is to fit $\theta^T \phi(x)$ with θ now a vector in \mathbb{R}^P , with the new update rule:

$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)})$$

In original update rule, we need to update every entry of θ and store it, but with kernel trick, we will not need to store θ explicitly due to the following fact:

$$\theta = \sum_{i=1}^n \beta_i \phi(x^{(i)})$$

[θ can be represented by linear combination of $(\phi(x^{(1)}), \dots, \phi(x^{(n)}))$]

and with induction we can obtain update rule for β_i :

$$\theta = \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)})$$

$$= \sum_{i=1}^n (\beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)}))) \phi(x^{(i)}) \Rightarrow \beta_i := \beta_i + \alpha y^{(i)} - \alpha \theta^T \phi(x^{(i)})$$

Moreover as we substitute θ : $\beta_i := \beta_i + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j \phi(x^{(j)})^T \phi(x^{(i)})) = \beta_i + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}))$

(Kernel trick become interesting as • We can pre-compute $K(x^{(i)}, x^{(j)})$ before loop starts Briefly, ... n).

• for some feature maps, computing $K(x^{(i)}, x^{(j)})$ is efficient.

without the necessity to explicitly refer to $\phi(x^{(i)})$

Algorithm (LMS with kernels)

1. Compute all $K(x^{(i)}, x^{(j)}) \triangleq \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ and obtain, $[K]_{mn}$ with $K_{ij} = K(x^{(i)}, x^{(j)})$

2. Loop: $\vec{\beta} := \vec{\beta} + \alpha (\vec{y} - K\vec{\beta})$

When making predictions, $\theta^T \phi(x) = \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x) = \sum_{i=1}^n \beta_i K(x^{(i)}, x)$

consider $K(\cdot, \cdot)$ defined by $K(x, z) = (x^T z + c)^k = \sum_{i,j=1}^n (x_i x_j)(z_i z_j) + \sum_{i=1}^d (\sqrt{c} x_i)(\sqrt{c} z_i) + c^2$

the corresponding feature mapping $\phi(x)$ contains all feature combination up to order 2

[Example for $d=3$, $\phi(x) = [x_1 x_1, x_1 x_2, x_1 x_3, x_2 x_1, x_2 x_2, x_2 x_3, x_3 x_1, x_3 x_2, x_3 x_3, \sqrt{c} x_1, \sqrt{c} x_2, \sqrt{c} x_3, c]^T$]

More broadly, kernel $K(x, z) = (x^T z + c)^k$ corresponds to an C_{k+d}^k feature space, corresponding to all monomials of form $x_{i_1} x_{i_2} \cdots x_{i_k}$ that are up to order k , notice in this $O(d^k)$ -dimensional

space, computing $K(x, z)$ still only takes $O(d)$ time

[Notice $C_{k+d}^k \approx (k+d)^k \rightarrow e$ when $k \rightarrow \infty$, think this could give an intuition of why

the Gaussian kernel $K(x, z) = \exp(-\frac{\|x-z\|^2}{2\sigma^2})$ can be considered as infinite dimensional]

* Necessary conditions for valid kernels

let $K \in \mathbb{R}^{n \times n}$ be defined s.t. $K_{ij} = K(x^{(i)}, x^{(j)})$, then the following theorem gives sufficient condition for K to be a valid kernel.

Theorem (Mercer): K is a valid (Mercer) kernel \Leftrightarrow for $\forall \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$, the corresponding kernel matrix is symmetric positive semi-definite.

Proof (\Rightarrow): For any $z \in \mathbb{R}^d$, we have $z^T K z = \sum_i \sum_j z_i k_{ij} z_j = \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(i)}) z_j$
 $= \sum_i \sum_j z_i \sum_k \phi_k(x^{(i)}) \phi_k(x^{(i)}) z_j = \sum_i \sum_j \sum_k z_i \phi_k(x^{(i)}) \phi_k(x^{(i)}) z_j = \sum_k \left(\sum_i z_i \phi_k(x^{(i)}) \right)^2 \geq 0$

L_1 -norm soft margin SVM

To make SVM algorithm work for non-linearly separable datasets as well as less susceptible to outliers, we can reformulate the optimization as:

$$\begin{aligned} & \min_{w, b, \gamma} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \gamma_i \\ & \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - \gamma_i, \quad i=1, 2, \dots, n \quad [\text{relax boundary and add penalty}] \\ & \quad \gamma_i \geq 0, \quad i=1, 2, \dots, n \end{aligned}$$

The dual form of problem then becomes:

$$\begin{aligned} & \max_{\alpha} W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ & \text{s.t. } 0 \leq \alpha_i \leq C \\ & \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0, \quad i=1, 2, \dots, n \end{aligned}$$

and prediction formulas remain the same, the KKT dual-complementarity conditions in this case are

$$\begin{cases} \alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1 \\ \alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1 \\ 0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1 \end{cases}$$

Cross Validation

simple (hold-out) cross validation

1. randomly split S into S_{train} (say 70% of data), S_{dev} (S_{cv} , the remaining 30%) [$(7-3)/(6-2-2)$]

2. train each model (M_i , where $M = \{M_1, M_2, \dots, M_n\}$ is a finite set of models we are selecting among)

on S_{train} to get some hypothesis h_i .

[choose $S_{\text{train}}/S_{\text{test}}$ in a way that let meaningful comparison happen between models]

3. select the h_i that has the smallest error $\hat{\epsilon}_{S_{\text{cv}}}^{(h_i)}$ on the development set S_{dev} .

Optionally, step 3 can be carried out as selecting M_i according to $\arg \min \{\hat{\epsilon}_{S_{\text{cv}}}^{(h_i)}\}$ and then retrain

M_i on the *ENTIRE* training set again.

[exceptions being algorithms sensitive to perturbations of initial data]

k -fold cross validation

[on relatively small training set]

1. randomly split S into k disjoint subsets each containing m/k training samples.

2. for each model M_i , do

For $j=1, 2, \dots, k$, train the model M_i on $S_1 \cup S_2 \cup \dots \cup S_{j-1} \cup S_{j+1} \cup \dots \cup S_k$
 to get hypothesis h_{ij} .

Test the hypothesis h_{ij} on S_j and calculate $\widehat{\mathbb{E}}_{S_j}(h_{ij})$, and let error of M_i being the average of $\widehat{\mathbb{E}}_S(h_{ij})$'s

3. Pick the Model M_i with $i = \operatorname{argmin} \frac{1}{K} \sum_{j=1}^K \widehat{\mathbb{E}}_{S_j}(h_{ij})$.

$K=10$ is a typical choice of K -fold cross validation, when data is really scarce, consider

* Leave-one cross validation with $K=m$ (the size of training set).

[Note:] Here we denote the training model set M as finite, but if we were to select from an infinite model set (e.g. choose bandwidth parameter $T \in \mathbb{R}^+$), we can 1. discretize $T \in \mathbb{R}^+$ to form finite selection space; 2. perform search over infinite model classes (optimization space $= \infty$).

[Note:] Cross validation can not only be used on selection of the model, but evaluation as well.

• Feature Selection

Apply feature selection algorithm to reduce the number of features.

* Algorithm Forward Search (Wrapper model feature selection)

1. Initialize $\mathcal{F} = \emptyset$

A: given the number of features each iteration, the model is decided, both are not decided?

2. Repeat {

(a) For $i=1, 2, \dots, d$ (d is the number of features), if $i \notin \mathcal{F}$, let $\mathcal{F}_i = \mathcal{F} \cup \{i\}$, train the learning model using only the features in \mathcal{F}_i and estimate the generalization error [can be done via cross-validation]

(b) Set $\mathcal{F} := \mathcal{F}_{\operatorname{argmin}} \widehat{\mathbb{E}}_S(h_i)$, which is the set that performs best in step (a).

} Select the best feature subset that was evaluated during the entire search procedure.

[Main idea: Adding feature greedily, and the outer loop can terminate with a pre-set $|F|$.]

* Filter Feature Selection

1. For discrete-valued features x_i , choose the informative score $S(i)$ to be mutual information:

$$S(i) := MI(x_i; y) = \sum_{x_i \in \mathcal{X}_i} \sum_{y \in \mathcal{Y}} p(x_i, y) \log \frac{p(x_i, y)}{p(x_i)p(y)} \quad [\text{extrapolation to finite set is trivial}]$$

2. More generally, notice mutual information can be expressed as Kullback-Leibler(KL) divergence:

$$MI(x_i; y) = KL(p(x_i, y) || p(x_i)p(y))$$

$$\text{where } KL(p||q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} = - \sum_{x \in X} p(x) \log \frac{q(x)}{p(x)}$$

in case where $q(x) \neq p(x)$
is unsuitable for numerator

For distributions p and q of a continuous random variable, relative entropy is defined as

$$KL(p||q) = \int_{-\infty}^{+\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx.$$

in terms of Bayesian inference, $KL(P||Q)$ is a measure of information gained by revising prior probability distribution Q to posterior probability distribution P (i.e. the info lost when Q is used to approximate P).

Note: the Neyman-Pearson Lemma states that the most powerful way to distinguish between distributions P and Q based on observation Y (drawn from one of them) is through $\log \frac{P(Y)}{Q(Y)}$,

KL divergence appears to be the expectation if Y is actually drawn from P , namely, $E[\log \frac{P(Y)}{Q(Y)}]$.

Logic chain: If x_i and y is independent, then KL divergence between $p(x_i|y)$ and $p(x_i)p(y)$ will be zero, thus $S(x_i)$ will be small, so x_i will be unlikely selected.

3. After ranking the scores, do a cross-validation to decide on the number of features to select.

Bayesian Statistics

Frequentist statistics view θ as being constant valued but unknown, and it's our job to come up with statistical procedures (e.g. maximum likelihood) to estimate θ .

Bayesian view θ as being a random variable, and our "prior belief" of θ is reflected on our choice of θ prior distribution $p(\theta)$. When given a data set $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$, we can compute the posterior distribution as:

$$p(\theta|S) = \frac{\prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta) p(\theta)}{\int_{\theta} \prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta) p(\theta) d\theta}$$

$\rightarrow p(\theta|S)$ depend on the model using

computationally expensive when integrating over high-dimensional θ .

When given a new test sample x , compute $p(y|x, S) = \int_{\theta} p(y|x, \theta) p(\theta|S) d\theta$.

and $E(y|x, S) = \int_y y p(y|x, S) dy$. [sum if y is discrete]. (maximum a posteriori)

Instead of computing for posterior distribution, we can do an approximation, epitome being MAP:

$$\theta_{MAP} = \arg \max_{\theta} \prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta) p(\theta) \quad [\text{single-point estimation}]$$

Common choice for prior $p(\theta)$ is to assume that $\theta \sim \mathcal{N}(0, \tau^2 I)$, using this choice, $\|\theta_{MAP}\|_2 < \|\theta_{MLE}\|_2$.

In practice, this cause Bayesian MAP being less susceptible to overfitting than ML estimate.

[due to regularization, and why Bayesian logistic regression is effective for text classification]

Bias - Variance Analysis

Say our goal is to construct an estimator for the unknown parameter θ^* given observed data set $\{(\mathbf{x}^{(i)}, y^{(i)})\}$.

The distribution of our model output $\hat{\theta}_n$ (random variable given noise typically in label \vec{y}) is called

sampling distribution, and bias and variance are first and second moments of its sampling distribution, namely $\text{Bias}(\hat{\theta}_n) \equiv E[\hat{\theta}_n - \theta^*]$, $\text{Var}(\hat{\theta}_n) \equiv \text{Cov}(\hat{\theta}_n)$

Bias and Variance of an estimator are not necessarily directly related, under squared error, they can be

related as: $\text{MSE}(\hat{\theta}_n) = E[\|\hat{\theta}_n - \theta^*\|^2] =$

$$= E[\|\hat{\theta}_n - E[\hat{\theta}_n] + E[\hat{\theta}_n] - \theta^*\|^2]$$

$$= E[\|E[\hat{\theta}_n] - \theta^*\|^2 + \|E[\hat{\theta}_n] - \theta^*\|^2 + 2(\hat{\theta}_n - E[\hat{\theta}_n])^T(E[\hat{\theta}_n] - \theta^*)]$$

$$= \|E[\hat{\theta}_n] - \theta^*\|^2 + E[\text{tr}[(\hat{\theta}_n - E[\hat{\theta}_n])(\hat{\theta}_n - E[\hat{\theta}_n])^T]]$$

$$= \text{tr}[\text{Var}(\hat{\theta}_n)] + \|\text{Bias}(\hat{\theta}_n)\|^2.$$

[quite often the case: techniques employed to reduce variance results in increased bias and vice versa]

Bias-Variance Tradeoff in linear Regression with L_2 -regularization.

[L_2 -norm]

In linear regression, when taking a probabilistic view, we assume $y^{(i)} = \theta^T \mathbf{x}^{(i)} + \epsilon^{(i)}$,

with each $\epsilon^{(i)} \sim N(0, \sigma^2)$ (i.i.d.). Then L_2 -regularization led to minimizing the cost function

$$J(\theta) = \frac{\lambda}{2} \|\theta\|_2^2 + \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2. \quad [\text{Regularization parameter } \lambda > 0]$$

and enjoys a closed-form solution $\hat{\theta}_n = \arg \min_{\theta \in \mathbb{R}^d} J(\theta) = \arg \min_{\theta \in \mathbb{R}^d} \left[\frac{\lambda}{2} \|\theta\|_2^2 + \frac{1}{2} \|\mathbf{x} \theta - \vec{y}\|_2^2 \right]$

$$\hat{\theta}_n \text{ s.t. } \nabla_{\theta} J(\hat{\theta}_n) = 0 \Rightarrow (\lambda I + \mathbf{x}^T \mathbf{x}) \hat{\theta}_n = \mathbf{x}^T \vec{y} \Rightarrow \hat{\theta}_n = (\lambda I + \mathbf{x}^T \mathbf{x})^{-1} (\mathbf{x}^T \vec{y})$$

Consider eigendecomposition of PSD matrix $\mathbf{x}^T \mathbf{x}$: $\mathbf{x}^T \mathbf{x} = U \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{bmatrix} U^T$ with $U^T U = I_n$

even though \mathbf{x} (hence $\mathbf{x}^T \mathbf{x}$) is not full rank, $(\mathbf{x}^T \mathbf{x} + \lambda I)$ is always symmetric and PD.

$$\Rightarrow (\mathbf{x}^T \mathbf{x} + \lambda I)^{-1} = U \begin{bmatrix} \frac{1}{\sigma_1^2 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sigma_n^2 + \lambda} \end{bmatrix} U^T \Rightarrow \hat{\theta}_n \text{ always exists and is unique.}$$

* Bias-Variance analysis: $\hat{\theta}_n = (\mathbf{x}^T \mathbf{x} + \lambda I)^{-1} \mathbf{x}^T (\mathbf{x} \theta^* + \vec{\epsilon}) = (\mathbf{x}^T \mathbf{x} + \lambda I)^{-1} \mathbf{x}^T \mathbf{x} \theta^* + (\mathbf{x}^T \mathbf{x} + \lambda I)^{-1} \mathbf{x}^T \vec{\epsilon}$

$$E[\hat{\theta}_n] = U \begin{bmatrix} \frac{1}{\sigma_1^2 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sigma_n^2 + \lambda} \end{bmatrix} U^T \mathbf{x}^T \theta^* = U \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{\sigma_n^2}{\sigma_n^2 + \lambda} \end{bmatrix} U^T \theta^*$$

$$E[\vec{\epsilon}] = 0 \quad \text{so} \quad \hat{\theta}_n = U \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{\sigma_n^2}{\sigma_n^2 + \lambda} \end{bmatrix} U^T \theta^*$$

Observations: 1. when $\lambda = 0$, then $E[\hat{\theta}_n] = \theta^* \Rightarrow$ unbiased but biggest variance.

2. the larger λ , the smaller the eigenvalues will be, shrink towards 0. \Rightarrow more bias

\Rightarrow the estimator $\hat{\theta}_n$ of L_2 -regularization Linear Regression is Biased.

$$\begin{aligned}
 \text{Cov}[\hat{\theta}_n] &= \text{Cov}\left[(\hat{x}^T x + \lambda I)^{-1} \hat{x}^T \hat{\epsilon} + (\hat{x}^T x + \lambda I)^{-1} \hat{x}^T \hat{\epsilon}\right] \\
 &= \text{Cov}\left[(\hat{x}^T x + \lambda I)^{-1} \hat{x}^T \hat{\epsilon}\right] \quad [\text{if } \hat{\epsilon} \sim N(0, \tau^2 I) \text{ then } A \hat{\epsilon} \sim N(0, A(\tau^2 I)A^T)] \\
 &= (\hat{x}^T x + \lambda I)^{-1} \hat{x}^T \text{Cov}[\hat{\epsilon}] x (\hat{x}^T x + \lambda I)^{-1} \\
 &= \tau^2 (\hat{x}^T x + \lambda I)^{-1} \hat{x}^T x (\hat{x}^T x + \lambda I)^{-1} \\
 &= \tau^2 U \begin{bmatrix} \frac{\sigma_1^2}{(\sigma_1^2 + \lambda)} & \cdots & 0 \\ 0 & \cdots & \frac{\sigma_n^2}{(\sigma_n^2 + \lambda)} \end{bmatrix} U^T
 \end{aligned}$$

Observations: 1. if larger λ , then smaller $\text{tr}(\text{Cov}[\hat{\theta}_n]) \Rightarrow$ less variance (but more bias)

Bias and Variance in Prediction

Given a training set $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$, under the assumption $\bar{y} = f(x) + \epsilon$ where $E[\epsilon] = 0, \text{Var}[\epsilon] = \tau^2$.

Further, we define the true f as $f(x) \equiv E[y|x=x]$. expected square error,

* Construct \hat{f}_n to mimic f on unseen examples \Rightarrow small generalization error for \hat{f}_n .

[Note]: \hat{f}_n is random due to the randomness embedded in $\epsilon^{(i)}$'s.

Consider a new unseen example pair (y^*, x^*) , we obtain:

$$\begin{aligned}
 \text{MSE}[\hat{f}_n] &= E[(y^* - \hat{f}_n(x^*))^2] = E[(f(x^*) - \hat{f}_n(x^*))^2 + \epsilon^2 - 2(f(x^*) - \hat{f}_n(x^*))] \\
 &= \tau^2 + E(f(x^*) - \hat{f}_n(x^*))^2 + \text{Var}(f(x^*) - \hat{f}_n(x^*)) = \tau^2 + \text{Bias}^2 + \text{Variance}.
 \end{aligned}$$

* relation with bias and variance in inference -

$$\text{Bias}[\hat{f}_n] = E[\hat{f}_n(x) - f(x)] = E[\hat{\theta}_n x - \theta^T x] = E[\hat{\theta}_n - \theta]^T x = \text{Bias}(\hat{\theta}_n)^T x.$$

and similarly, $\text{Variance}(\hat{f}_n) = x^T [\text{Var}(\hat{\theta}_n)] x$.

the irreducible error only exists in test setting, since it is an artifact of noise in test sample.

[noise in the training data contributes to Variance term, noise in test sample manifests as irreducible error].

+ in inference setting, observed bias-variance tradeoff

+ in prediction setting, choose λ to be the value that minimizes the squared error in cross-validation. (L2 regularization, LG)

* relation with overfitting and underfitting .

+ overfitting relates to having high variance $\xrightarrow{\text{fight overfit}}$ { increase regularization
increase data size
decrease # features
use a smaller model. }

+ underfitting relates to having high bias $\xrightarrow{\text{fit underfit}}$ { decrease regularization
use more features
use a larger model. }

+ training error can be treated as the amount of bias in model / estimator. [high bias \rightarrow underfitting]

+ gap between cross-validation error and training error can be treated as Variance.

We should always analyse the model performance by looking at training error and cross-validation error simultaneously. (then take steps to address either Bias or Variance purposefully)

Note: Steps taken to fight overfitting (i.e. high Variance) generally do not help fight underfitting, it is futile, for example, to obtain more data (fight high Variance) when the training error is high (symptom of high Bias), vice versa.

Decision Trees

select leaf node, feature, threshold

ONE at a time.

Generate approximate solution of region via greedy, top-down, recursive partitioning.

take the loss function L start at original input region X

as splitting heuristic and split into child regions by

THRESHOLDING on SINGLE feature.

given a parent region R_p , feature index j , and threshold t_{CR} , we obtain child regions R_1, R_2 as:

$$R_1 = \{x | x_j < t, x \in R_p\}, \quad R_2 = \{x | x_j \geq t, x \in R_p\}$$

Within greedy partitioning framework, we want to select (leaf, feature, threshold) that maximize the decrease in loss: $L(R_p) - (|R_1|L(R_1) + |R_2|L(R_2)) / (|R_1| + |R_2|)$ [child region decrease the loss of R_p]

specifically for classification problem, misclassification loss on R : $L_{mis}(R) = 1 - \max_c(\hat{p}_c)$. (1)

where \hat{p}_c denotes the samples $\in R$ that are of class c . [the proportion of samples that would be missed if predicting the majority class for region R]

We notice that misclassification loss is not very sensitive to changes in class probabilities, propose a more sensitive loss (cross-entropy loss): $L_{cross}(R) = - \sum_c \hat{p}_c \log \hat{p}_c$ [information gained through splitting]

[number of bits needed to specify the outcome (class) given the distribution is known]

since the cross-entropy loss is strictly concave, as long as $\hat{p}_1 \neq \hat{p}_2$ (in a binary classification problem),

$L(R_p) > (|R_1|L(R_1) + |R_2|L(R_2)) / (|R_1| + |R_2|)$, which make cross-entropy loss more sensitive/informative.

For regression settings, we can directly use squared loss to select splits:

$$\text{final } L_{loss} = L_{\text{squared}}(R) = \sum_{i \in R} (y_i - \hat{y})^2 / |R|$$

and prediction for region R now becomes: $\hat{y} = \sum_{i \in R} y_i / |R|$. \uparrow plugin

• Regularization

Decision tree is high varianced and low biased, thus we turn to techniques of regularization:

* Minimum Leaf Size. Stop splitting R if cardinality falls below a fixed threshold.

* Maximum Depth. Stop splitting R if more than a fixed threshold of splits were taken to reach R.

* Maximum Number of Nodes. Stop if tree has more than a fixed threshold of leaf nodes.

* Pruning. * Fully Grown Out * the tree and then pruning away nodes that minimally decrease loss.

Notice it's not a good idea to prune along the growing process since single-feature selection would miss higher order interactions if constrained on loss decrease.

• Ensemble Learning

Some of the primary benefits of decision trees are: categorical variable support, fast, interpretable.

While disadvantages include: * High variance *(imagine a separate region for every data point)

* Poor Additive Modeling *(compare with simple linear model)

leading to low predictive accuracy of individual decision tree.

→ * Bagging (Bootstrap Aggregation) ← [decrease Variance]

Main idea: assume training set is the true population (P), let $\text{card}(S) = n$, then sample $\downarrow \vec{x}$ bootstrap samples from S , and do it with replacement (always holds $|S| = n$) \uparrow unknown

Bootstrap samples $\vec{z}_1, \vec{z}_2, \dots, \vec{z}_m$, train model G_i on \vec{z}_i then $G(\vec{x}) = \sum_{i=1}^m G_i(\vec{x}) / m$. second term goes down as $m \uparrow$.

Bootstrap drives down correlation between X 's, since $\text{Var}(\vec{x}) = p\sigma^2 + \frac{1-p}{m}\sigma^2 \rightarrow$ sample more models until the second term stops changing.

Note: The tradeoff for reducing variance is that p decrease as # samples increase
models get less complicated so bias goes up. but has a lower bound since samples drawn can't be completely uncorrelated.

→ * Random Forest ←

1. Introduce more randomization into each decision tree to further decorrelate the samples.

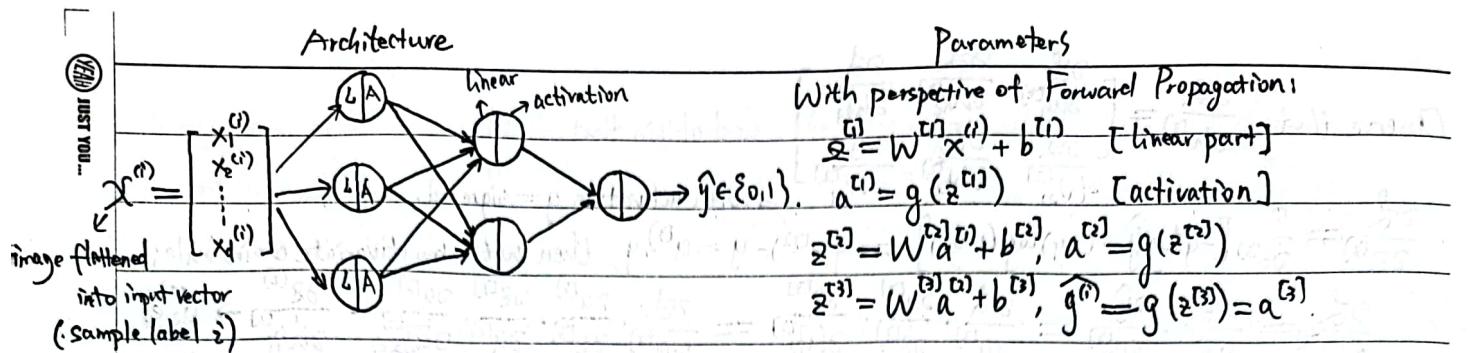
2. At each split, consider only a fraction of total features. (typically \sqrt{p} features per split)

randomly select features sampling with replacement

Neural Network

Hidden Layer: do not have the ground truth/training value for hidden units, the network itself is responsible for figuring out intermediate structures. Let $a_j^{(i)}$ denote the output from the j^{th} hidden unit in the i^{th} layer. ($i \geq 1, i=0$ denoting input layer).

daallen°



In neural networks, with purpose of optimizing training speed, *vectorization* is performed in substitution of for loops, here we denote

$$\begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{bmatrix} = \underbrace{\begin{bmatrix} -W_1^{(1)} \\ -W_2^{(1)} \\ -W_3^{(1)} \end{bmatrix}}_{\underline{W}^{(1)} \in \mathbb{R}^{3 \times d}} \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ \vdots \\ x_d^{(1)} \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix},$$

with a total of 3d+4 parameters

(With row being set of parameter for neuron)

• Initializing Parameters

First, note that we can't initialize all parameters to zero, this would cause output of first layer to be zero regardless of input. Moreover, the output of network would always be 0.5 regardless of input. One solution is to randomly initialize parameters to small values (e.g. $\mathcal{N}(0, 0.1)$).

In practice, there is something better than random initialization, namely Xavier/He initialization which initializes weights: $W^{(l)} \sim \mathcal{N}(0, \sqrt{\frac{2}{n^{(l)} + n^{(l+1)}}})$, where $n^{(l)}$ is the number of neurons in layer l . [Intuition: consider variance (input to layer l) to be $\sigma^{(in)}$ and variance (output of layer l) to be $\sigma^{(out)}$, Xavier/He initialization encourages $\sigma^{(in)}$ being SIMILAR to $\sigma^{(out)}$]. Also, note that we can't initialize all parameters to be the same non-zero value, which would cause each element of $a^{(l)}$ to be the same, as a result, when computing gradient, all neurons in a layer will be equally responsible for anything contributed to final loss.

then each neuron will be updated the same way and learn the exact same thing. [symmetry]

• Optimization

logistic regression for binary classification

Using the Logistic Loss $\mathcal{L}(\hat{y}, y) = -[(1-y)\log(1-\hat{y}) + y\log\hat{y}]$, we can perform SGD to optimize parameters in each layer l :

$$W^{(l)} := W^{(l)} - \alpha \frac{\partial \mathcal{L}}{\partial W^{(l)}}(1), \quad b^{(l)} := b^{(l)} - \alpha \frac{\partial \mathcal{L}}{\partial b^{(l)}}(1)$$

the derivations we need to compute are $(\frac{\partial \mathcal{L}}{\partial W^{(l)}}, \frac{\partial \mathcal{L}}{\partial b^{(l)}}) (\frac{\partial \mathcal{L}}{\partial w^{(l)}}, \frac{\partial \mathcal{L}}{\partial b^{(l)}})$.

We will be calculating $\frac{\partial \mathcal{L}}{\partial w^{(1)}}$ as an example.

Note: Loss function selection: ① regression $\mathcal{L}(\hat{y}, y) = \frac{1}{2} \sum_i (y^{(i)} - \hat{y}^{(i)})^2$.

② softmax regression - cross entropy loss $\mathcal{L}(\hat{y}, y) = - \sum_i 1 \cdot y^{(i)} \log \hat{y}^{(i)}$

serve that, $\frac{\partial \mathcal{L}}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_{11}^{(2)}} & \frac{\partial \mathcal{L}}{\partial w_{12}^{(2)}} & \frac{\partial \mathcal{L}}{\partial w_{13}^{(2)}} \\ \frac{\partial \mathcal{L}}{\partial w_{21}^{(2)}} & \frac{\partial \mathcal{L}}{\partial w_{22}^{(2)}} & \frac{\partial \mathcal{L}}{\partial w_{23}^{(2)}} \end{bmatrix}$, and obtain that

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^{(2)}} = \frac{\partial}{\partial g^{(2)}} [-y \log \hat{y} - (1-y) \log (1-\hat{y})] = g(z^{(2)}) - y = a^{(2)}_j - y, \text{ then with multivariate chain rule:}$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^{(2)}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{ij}^{(2)}} = \frac{\partial \mathcal{L}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w_{ij}^{(2)}} = \underbrace{\frac{\partial \mathcal{L}}{\partial a^{(2)}}}_{\substack{(a^{(2)} - y) \\ R^{1 \times 1}}} \cdot \underbrace{\frac{\partial a^{(2)}}{\partial z^{(2)}}}_{\substack{(W^{(2)}) \\ R^{1 \times 2}}} \cdot \underbrace{\frac{\partial z^{(2)}}{\partial w_{ij}^{(2)}}}_{\substack{g(z^{(2)}) (1-g(z^{(2)})) \\ R^{2 \times 2}}} \rightarrow a_j^{(2)} e_i^{(2)}$$

$$= \underbrace{(a^{(2)} - y) W^{(2)} \circ g'(z^{(2)})}_{R^{1 \times 2}} \underbrace{a_j^{(2)} e_i^{(2)}}_{R^{2 \times 1}}$$

$$= [(a^{(2)} - y) W^{(2)} \circ g'(z^{(2)})] a_j^{(2)} \rightarrow \frac{\partial \mathcal{L}}{\partial w^{(2)}} = [(a^{(2)} - y) W^{(2)} \circ g'(z^{(2)})] a^{(2)T}$$

Notice when calculating $\frac{\partial \mathcal{L}}{\partial w^{(2)}} = \frac{\partial \mathcal{L}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(2)}} = \frac{\partial \mathcal{L}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(2)}} = (a^{(2)} - y) a^{(2)T}$, the first term has been calculated in the derivation of $\frac{\partial \mathcal{L}}{\partial w^{(2)}}$; and $\frac{\partial \mathcal{L}}{\partial w^{(2)}} = \frac{\partial \mathcal{L}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial w^{(2)}} = \frac{\partial \mathcal{L}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial w^{(2)}}$ calculated in $\frac{\partial \mathcal{L}}{\partial w^{(2)}}$

The idea is, in order for errors to back propagate, we need to go through variables that are connected to each other and apply the multivariate chain rule of calculus.

Notice: Notice when expressing $\frac{\partial \mathcal{L}}{\partial w^{(2)}}$ we end up with $\frac{\partial \mathcal{L}}{\partial w^{(2)}} = W^{(2)T} \circ a^{(2)} (1-a^{(2)}) (a^{(2)} - y) a^{(2)T}$, where $a^{(2)}, a^{(2)T}$ are values there were computed during forward propagation; it's important to store all values that are computed during forward propagation for efficient backward update.

When optimizing variables, we usually won't choose stochastic gradient descent since loss for a single example might be noisy, while $J = \frac{1}{n} \sum_{i=1}^n \mathcal{L}^{(i)}$ gives more accurate gradients but is computationally expensive. In practice, prioritize using mini-batch gradient descent, with $J_{mb} = \frac{1}{B} \sum_{i=1}^B \mathcal{L}^{(i)}$, where B is the number of examples in the mini-batch.

Momentum: (for training phase optimization) Consider mini-batch stochastic gradient, for any layer i , the update rule becomes:

$$\nabla_{w^{(i)}} = \beta \nabla_{w^{(i)}} + (1-\beta) \frac{\partial \mathcal{L}}{\partial w^{(i)}}$$

$$w^{(i)} = w^{(i)} - \nabla_{w^{(i)}}$$

the weight update now depends on two terms: the cost \mathcal{L} at this update step and velocity $\nabla_{w^{(i)}}$, with relative importance controlled by β . [velocity $\nabla_{w^{(i)}}$ will keep track of gradient over time].

Choice of Activation Functions:

One way to improve the performance of neural network is to use different activation functions.

Usually choose to use same activation function within the same layer, with some common choices being:

Sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$ ReLU (Rectified Linear Unit) $r(x) = \max(x, 0)$ tanh $(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
 $\sigma'(x) = \sigma(x)(1-\sigma(x))$. $r'(x) = \begin{cases} 1 & x > 0 \\ 0 & \text{otherwise} \end{cases}$. $\tanh'(x) = 1 - \tanh^2(x)$.

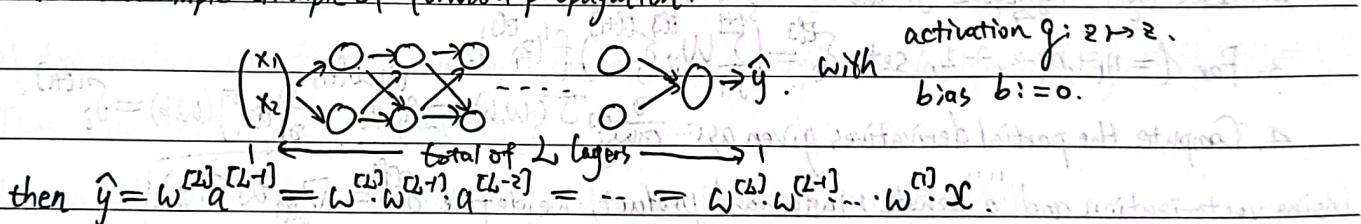
Reasons to choose these activations come along with different tasks. When doing classification work with labeling being {0, 1}, then sigmoid is a common choice; In reinforcement learning when there is necessity for negative reward, tanh would be a reasonable choice; When performing regression tasks,

ReLU might appear to be an ideal choice. The problem with sigmoid is that gradients are small when x becomes big/small, namely gradient vanishes at two sides, making it hard to train network at the early stage since gradient drops too slowly. In this case, having a ReLU layer as the first hidden layer might be a good adjustment.

[Note]: one initialization technique to avoid input being too big/small is perform $\tilde{x} = \frac{x - \mu_{\text{train}}}{\sigma_{\text{train}}}$, what this causes is the gradient of loss function (contour) always points to the center.

• Vanishing/Exploding Gradients.

Derive a simple example of forward propagation:



if we assume $w^{[l]} = \text{diag}(1+\varepsilon, 1+\varepsilon)$, with $\varepsilon > \frac{1}{2} \Rightarrow \hat{y} = \text{diag}((1+\varepsilon)^L, (1+\varepsilon)^L) x$, if L is large(enough),

then \hat{y} would explode, similarly $\varepsilon < -\frac{1}{2}$, \hat{y} will be extremely small if layer numbers are large.

Tackling: we would prefer weights to be as close to 1 as possible, so a reasonable(not perfect) way is to initialize input and weights properly. With sigmoid activation on layer l , set $w^{[l]} = \text{shape}^{[l]} \times \sqrt{\frac{1}{n^{[l-1]}}}$. (With ReLU as activation, $w^{[l]} = \text{randn}(\text{shape}^{[l]}) \times \sqrt{\frac{2}{n^{[l-1]}}}$ randomly set).

With tanh as activation, Xavier Initialization proposes $w^{[l]} \sim N(0, \sqrt{\frac{1}{n^{[l-1]}}})$, if we also consider backward gradient, then He initialization with $w^{[l]} \sim N(0, \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}})$ would be a good choice.

Regularization and Parameter Sharing

Pile all parameters in layer 1 to form a single matrix W . Add L_2 -norm regularization to the cost function:

$$J = J + \frac{\lambda}{2} \|W\|_2^2 = J + \frac{\lambda}{2} \sum_{ij} |W_{ij}|^2 = J + \frac{\lambda}{2} W^T W.$$

where λ is an arbitrary value with a larger value indicating more regularization, then the update rule becomes: $W = W - \alpha \frac{\partial J}{\partial W} - \alpha \frac{\lambda}{2} \frac{\partial W^T W}{\partial W} = (1-\alpha\lambda)W - \alpha \frac{\partial J}{\partial W}$. This indicates that with L_2 regularization, every update will include penalization ($\alpha\lambda W$) depending on W .

Recall in logistic regression, the linear part $z = \theta^T x$ demands θ_0 always look at the top left pixel (x_0), if all training sample appear at the right bottom, then θ_0 will never get trained. In order to let θ_0 gain more insights about image, we turn to convolutional neural networks, and let θ be a matrix (kernel), and activation becomes $a = \theta * x$. [Intuition: let θ "seen" all pixels of image].

Backpropagation

For each node i in layer l , let $\delta_i^{(l)}$ denote how much the node is responsible for errors in output.

Then backpropagation algorithm can be described as:

1. Perform a feedforward pass, compute activations $a^{(l)}$ for all layers ($l = 0, 1, 2, \dots, n_f$)

2. Let $\delta_i^{(n_f)} = \frac{\partial J}{\partial z_i^{(n_f)}} \frac{1}{2} \|g - h_{W,b}(x)\|^2 = (a_i^{(n_f)} - y_i) \cdot f'(z_i^{(n_f)})$.

3. For $l = n_f-1, n_f-2, \dots, 1$, set $\delta_i^{(l)} = \left(\sum_{j=1}^{n_{l+1}} W_{ji} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$.

4. Compute the partial derivatives given as: $\frac{\partial J}{\partial w_{ij}^{(l)}} = a_j^{(l)} \delta_i^{(l)}, \frac{\partial J}{\partial b_i^{(l)}} = \delta_i^{(l)}$.

Using vectorization and " \circ " denote Hadamard product, we define $\delta^{(l)} = \nabla_{z^{(l)}} \mathcal{L}(\hat{y}, y)$.

For the output layer N we have $\delta^{(N)} = \nabla_{z^{(N)}} \mathcal{L}(\hat{y}, y) = (a^{(N)} - y) \circ g'(z^{(N)})$.

For $l = N-1, N-2, \dots, 1$, we have $\delta^{(l)} = (W^{(l+1)T} \delta^{(l+1)}) \circ g'(z^{(l)})$, and the gradients become:

$$\nabla_{w_{ij}^{(l)}} J(w, b) = \delta^{(l+1)} a^{(l)T}, \nabla_{b_i^{(l)}} J(w, b) = \delta^{(l+1)} i$$

For output layer N , $\delta^{(N)} = \nabla_{z^{(N)}} \mathcal{L}(\hat{y}, y) = \underbrace{\nabla_{z^{(N)}} \mathcal{L}(\hat{y}, y) \circ g'(z^{(N)})}_{\text{elementwise derivative w.r.t. } z^{(N)}}$

For $l = N-1, N-2, \dots, 1$, $\delta^{(l)} = (W^{(l+1)T} \delta^{(l+1)}) \circ g'(z^{(l)})$, $\nabla_{w_{ij}^{(l)}} J(w, b) = \delta^{(l+1)} a^{(l-1)T}, \nabla_{b_i^{(l)}} J(w, b) = \delta^{(l+1)}$.

Learning Theory

To begin foray into learning theory, let's restrict attention to binary - classification in which labels are $y \in \{0, 1\}$. Assume training set being $S = \{(x^{(i)}, y^{(i)})\}_{i=1,2,\dots,n}$, where training samples $(x^{(i)}, y^{(i)})$ drawn independently identically distributed (iid) from probability distribution D .

Define training error (empirical risk) to be $\hat{\epsilon}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x^{(i)}) \neq y^{(i)}\}$, for a hypothesis h .

Define generalization error $\epsilon(h) = P_{(x,y) \sim D} (h(x) \neq y)$. [generate test example from D , probability of missing]

Note: Assumption of test example drawn from same distribution is one of the PAC assumptions, among which ① training and testing on same distribution ② independently drawn training examples were the most important. [probably approximately correct].

Consider setting of linear regression, let $h_\theta(x) = \mathbb{1}\{\theta^T x \geq 0\}$, with *empirical risk minimization* we select $\hat{\theta} = \arg \min_{\theta} \hat{\epsilon}(h_\theta)$.

- Define hypothesis class H to be set of all classifiers considered by the learning algorithm.

Then ERM can be considered to happen within H , thus $\hat{h} = \arg \min_{h \in H} \hat{\epsilon}(h)$.

- Consider the case of $|H|$ being finite, say H is set of k functions mapping X to $\{0, 1\}$.

Take any fixed $h_i \in H$, let random variable $Z = \mathbb{1}\{h_i(x) \neq y\}$ where $(x, y) \sim D$, then Z is an indicator of test sample misclassification, similarly define training sample misclassification $Z_j = \mathbb{1}\{h_i(x^{(j)}) \neq y^{(j)}\}$.

Thus $\hat{\epsilon}(h_i) = \frac{1}{n} \sum_{j=1}^n Z_j$, $\epsilon(h_i) = \mathbb{E}[Z]$.

- Lemma. (Union Bound) $P(A_1 \cup A_2 \cup A_3 \dots \cup A_k) \leq P(A_1) + P(A_2) + \dots + P(A_k)$.

- Lemma. (Hoeffding Inequality) Let Z_1, Z_2, \dots, Z_n be n iid random variables drawn from a Bernoulli(ϕ) distribution, let $\hat{\phi} = \frac{1}{n} \sum_{i=1}^n Z_i$, and let $\gamma > 0$ be fixed, then

$$P(|\hat{\phi} - \phi| > \gamma) \leq 2 \exp(-2\gamma^2 n).$$

(with Chernoff bound (Hoeffding)). We obtain $P(|\epsilon(h_i) - \hat{\epsilon}(h_i)| > \gamma) \leq 2 \exp(-2\gamma^2 n)$. (for one h_i)

Our goal now is to extrapolate from one particular h_i to all $h_i's \in H$. To do so, let's denote

A_i as the event $|\epsilon(h_i) - \hat{\epsilon}(h_i)| > \gamma$, then with Eq(1):

$$P(\exists h_i \in H, |\epsilon(h_i) - \hat{\epsilon}(h_i)| > \gamma) = P(A_1 \cup A_2 \cup A_3 \dots \cup A_k) \leq 2k \exp(-2\gamma^2 n)$$

$$\Rightarrow P(\forall h_i \in H, |\epsilon(h_i) - \hat{\epsilon}(h_i)| \leq \gamma) \geq 1 - 2k \exp(-2\gamma^2 n) \quad [\text{uniform convergence result}]$$

let $\delta = 2k \exp(-2\gamma^2 n)$, then with any two within $[n, \gamma, \delta]$ fixed, we can solve the bound

for the other. E.g. Given r and $\delta > 0$, how large must n be before we can guarantee that with probability at least $1-\delta$, training error will be within r for generalization error?

$$\Rightarrow \delta \geq 2k\exp(-2rn) \text{ for fixed } \delta, r, \text{ then } n \geq \frac{1}{2r^2} \log \frac{2k}{\delta} \quad [\text{sample complexity}]$$

Similarly, let n and δ be fixed, then with probability $1-\delta$, we have that for all $h \in H$,

$$|\hat{\epsilon}(h) - \epsilon(h)| \leq \sqrt{\frac{1}{2n} \log \frac{2k}{\delta}}.$$

Recall $\hat{h} = \arg \min_{h \in H} \hat{\epsilon}(h)$ is the output hypothesis that minimizes empirical error, define $h^* = \arg \min_{h \in H} \epsilon(h)$

the hypothesis that performs best on test sample, we observe a bound for $\hat{\epsilon}(\hat{h})$ using $\epsilon(h^*)$:

$$\hat{\epsilon}(\hat{h}) \leq \hat{\epsilon}(h^*) + r \leq \hat{\epsilon}(h^*) + \epsilon(h^*) + 2r.$$

- Theorem. Let $|H| = k$ and let any n, δ be fixed, with probability at least $1-\delta$, the inequality holds:

$$\hat{\epsilon}(\hat{h}) \leq \min_{h \in H} \epsilon(h) + 2\sqrt{\frac{1}{2n} \log \frac{2k}{\delta}}. \quad (2)$$

Using Ineq(2), we can have a interpretation of bias-variance tradeoff. Say we switched to a larger hypothesis set $H' \supseteq H$. Then first term in (2) will not increase, hence, by learning using a larger hypothesis set, "bias" will only decrease (not strictly). However, as $|H'|$ increases, second term in (2) also increase, this increase corresponds to "variance" increasing when using a larger hypothesis class.

- Corollary (sample complexity bound). Let $|H| = k$, and let any δ, r be fixed, then for $\hat{\epsilon}(\hat{h}) \leq \min_{h \in H} \epsilon(h) + 2r$ to hold with probability at least $1-\delta$, it suffices that

$$n \geq \frac{1}{2r^2} \log \frac{2k}{\delta} = O\left(\frac{1}{r^2} \log \frac{k}{\delta}\right) = O_{r, \delta}(k).$$

The Case of infinite H .

- Define. Vapnik-Chervonenkis dimension of hypothesis class H , written $VC(H)$ to be size of largest set that is shattered by H . Given a set $S = \{x^{(1)}, \dots, x^{(D)}\}$ of points $x^{(i)} \in \mathcal{X}$ (input space), we say that H shatters S if for any set of labels $\{y^{(1)}, y^{(2)}, \dots, y^{(D)}\}$ there exists some $h \in H$ s.t. $h(x^{(i)}) = y^{(i)}$ for $\forall i = 1, 2, \dots, D$. (Notice the set S is given in advance, not arbitrary.)

[Note:] In order to prove that $VC(H)$ is at least D , we need to show only there's at least one set of cardinality D that H can shatter.

- Theorem. Let H be given and $VC(H) = D$, then with probability at least $1-\delta$, we have that for

$$\text{all } h \in H, |\epsilon(h) - \hat{\epsilon}(h)| \leq O\left(\sqrt{\frac{D}{n} \log \frac{n}{D}} + \frac{1}{n} \log \frac{1}{\delta}\right) \Rightarrow \hat{\epsilon}(\hat{h}) \leq \epsilon(h^*) + O\left(\sqrt{\frac{D}{n} \log \frac{n}{D}} + \frac{1}{n} \log \frac{1}{\delta}\right).$$

In other words, if a hypothesis class has finite VC dimension, then uniform convergence occurs as n becomes larger.

-Corollary. For $|\hat{E}(h) - \bar{E}(h)| \leq \gamma$ to hold for $\forall h \in \mathcal{H}$. (hence $\hat{E}(h) \leq E(h^*) + 2\gamma$ with probability at least $1-\delta$) it suffices that $n = O_{\gamma, \delta}(D)$.

To conclude, for a given hypothesis class \mathcal{H} and for algorithm that tries to minimizes training error, # training samples needed to achieve generalization error close to that of optimal classifier is usually roughly linear in the number of parameters ($VC(\mathcal{H})$) of \mathcal{H} .

K-means Clustering

We now begin our foray into unsupervised learning problem. In clustering problem, we are given a training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ without labeling, and wish to group the data into cohesive "clusters".

Before writing out the algorithm for k-means clustering, we should decide the value of k . There are Information-Theoretic criteria like Akaike information criterion (AIC) to determine optimal value of k , but in practice, k can be decided using prior knowledge of properties of data set and implementation needs.

The k-means clustering algorithm:

1. Initialize cluster centroids $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$. (metric: randomly select k points from data set, say $\mathcal{X}_{init} = \{x^{(1)}, x^{(2)}, \dots, x^{(k)}\}$ and set $\mu_k := x^{(k)}$)

2. Repeat until convergence:

For every $i \in \{1, 2, \dots, m\}$, set $c^{(i)} := \arg \min_j \|x^{(i)} - \mu_j\|_2^2$ [color the point]

For each $j \in \{1, 2, \dots, k\}$, set $\mu_j := \frac{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\}}$.

[update μ_j to the center of points that share the j^{th} color].

In order to show the k-means algorithm converges, we will define the distortion function to be:

$$J(c, \mu) = \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|_2^2$$

Thus, J measures the sum of squared distances between each training example $x^{(i)}$ and the cluster centroid $\mu_{c(i)}$ to which it has been assigned. Let's take a digression to talk about the coordinate ascent algorithm.

Coordinate Ascent

Consider the unconstrained optimization problem $\max_{\alpha} W(\alpha_1, \dots, \alpha_n)$. Apart from gradient ascent and Newton's method, we introduce the coordinate ascent algorithm:

Loop until convergence : {

$$\text{For } i=1,2,\dots,n, \{ d_i := \arg \max_{d_i} W(d_1, \dots, d_{i-1}, \hat{d}_i, d_{i+1}, \dots, d_n) \}$$

Inside the loop, we hold all variables except for some d_i fixed, and reoptimize W w.r.t. just the parameter d_i ; the ordering can be decided using heuristics to see which variable expect to allow us make the largest increase in $W(\alpha)$.

Note: When W happens to be of such a form that "argmax" can be performed efficiently, coordinate ascent can be a fairly efficient algorithm.

It can be shown that k-means algorithm is coordinate ascent performed on $J(c, \mu)$:

$$\mu_j := \arg \min_{\mu_j} J(c, \mu) = \arg \min_{\mu_j} \sum_{i=1}^m \|x^{(i)} - \mu_j\|_2^2 \quad [\text{since } J \text{ is quadratic w.r.t. } \mu_j, \text{ it makes sense to employ coordinate descent}]$$
$$\text{Thus, } \nabla_{\mu_j} J(\mu) = 0 \Rightarrow \nabla_{\mu_j} \sum_{i=1}^m \mathbb{1}\{c^{(i)} = j\} (x^{(i)} - \mu_j)^T (x^{(i)} - \mu_j) = 0 \Rightarrow \mu_j := \frac{\sum_{i=1}^m \mathbb{1}\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{c^{(i)} = j\}}$$

[the update formula for $c^{(i)}$ can be obtained similarly when holding μ fixed and minimizing J w.r.t. $c^{(i)}$]

Thus, J must monotonically decrease. But notice J is non-convex, and so coordinate descent is not guaranteed to converge to global minimum. One common thing to do is run k-means many times, * using different random initial value for cluster centroids, then pick the clustering that gives the lowest distortion $J(c, \mu)$.

Expectation-Maximization Algorithm

Mixture of Gaussians.

Suppose we are given training set $\{x^{(1)}, \dots, x^{(m)}\}$ and wish to fit Gaussians to "clusters". Formally, we should specify a joint distribution $p(x^{(i)}, z^{(i)}) = p(x^{(i)} | z^{(i)}) p(z^{(i)})$, with $z^{(i)}$ being latent (unobserved) variables. Here, $z^{(i)} \sim \text{Multinomial}(\phi)$ (where $\phi \geq 0$, $\sum_j \phi_j = 1$ and parameter ϕ_j gives $p(z^{(i)} = j)$) and $x^{(i)} | z^{(i)} = j \sim \mathcal{N}(\mu_j, \Sigma_j)$. We let k denote the number of values that $z^{(i)}$ can take on (e.g. if $z^{(i)} \sim \text{Bernoulli}(\phi)$ then $k=2$ which is the case of GDA). Thus, our model posits that each $x^{(i)}$ was generated by ① randomly choosing $z^{(i)}$ from $\{1, 2, \dots, k\}$ and ② $x^{(i)}$ was drawn from one of k Gaussians depending on $z^{(i)}$.

Notice if $z^{(i)}$'s are unknown, then the maximization of

$$f(\phi, \mu, \Sigma) = \sum_{i=1}^m \log p(x^{(i)}; \phi, \mu, \Sigma) = \sum_{i=1}^m \log \sum_{j=1}^k p(x^{(i)} | z^{(i)}) p(z^{(i)} | \phi) \quad [\text{summing over } k \text{ possibilities}]$$

is relatively hard, but if like GDA we knew what the $z^{(i)}$'s were, then likelihood can be simplified as:

$$l(\phi, \mu, \Sigma) = \sum_{i=1}^m \sum_{j=1}^K 1\{z^{(i)}=j\} [\log p(x^{(i)}|z^{(i)}; \mu, \Sigma_j) + \log p(z^{(i)}=j)]$$

Maximizing this w.r.t. ϕ, μ and Σ gives the parameters:

$$\cdot \phi_j = \frac{1}{m} \sum_{i=1}^m 1\{z^{(i)}=j\}, \quad \cdot \mu_j = \frac{\sum_{i=1}^m 1\{z^{(i)}=j\} x^{(i)}}{\sum_{i=1}^m 1\{z^{(i)}=j\}}, \quad \cdot \Sigma_j = \frac{\sum_{i=1}^m 1\{z^{(i)}=j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m 1\{z^{(i)}=j\}}$$

EM algorithm comes in as we try to determine the value of $z^{(i)}$, the E-step in EM tries to "guess" the values of $z^{(i)}$ and pass it to the M-step. Iterative EM algorithm:

Repeat until convergence:

$$(E\text{-step}) \text{ For each } i, j, \text{ set } w_j^{(i)} := p(z^{(i)}=j|x^{(i)}; \phi, \mu, \Sigma) = \frac{p(x^{(i)}|z^{(i)}=j)p(z^{(i)}=j)}{\sum_{j=1}^K p(x^{(i)}|z^{(i)}=j)p(z^{(i)}=j)}$$

(M-step) Update the parameters:

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)}, \quad \mu_j := \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}, \quad \Sigma_j := \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}}$$

Basically, we are substituting $1\{z^{(i)}=j\}$ with $\mathbb{E}[1\{z^{(i)}=j\}] = w_j^{(i)}$ in the M-step. The values $w_j^{(i)}$ calculated in E-step represent "soft" guesses for values of $z^{(i)}$. [probability instead of direct assignment as we did in k-means]. Similar to k-means algorithm, EM is also susceptible to local optima, so it will be necessary to reinitialize at several different initial parameters and choose ϕ, μ, Σ that achieves $\max l(\phi, \mu, \Sigma)$.

• EM Convergence

Say we wish to fit a model $p(x, z)$ to the data $\{x^{(1)}, \dots, x^{(m)}\}$, where the likelihood is given by

$l(\theta) = \sum_{i=1}^m \log p(x^{(i)}; \theta) = \sum_{i=1}^m \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta)$, as we introduce latent variables, taking the joint distribution and marginalizing out these latent variables. [if $z^{(i)}$ is continuous, then substitute with integral]. Explicitly maximizing the estimate of θ may be hard given $z^{(i)}$ are unobserved. In this setting, EM algorithm proposes to repeatedly construct a lower-bound on $l(\theta)$ (E-step) and then optimizing this constructed lower bound. (M-step).

For each i , let $Q_i(z)$ be some distribution (left to be chosen) over the $z^{(i)}$'s, s.t. $\sum_z Q_i(z) = 1$.

$Q_i(z) \geq 0$, using Jensen's inequality we obtain:

$$\begin{aligned} \sum_i \log p(x^{(i)}; \theta) &= \sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) = \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \\ &\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log p(x^{(i)}, z^{(i)}; \theta) / Q_i(z^{(i)}). \end{aligned} \quad (1)$$

1. log being concave

2. Jensen: $f(\mathbb{E}X) \geq \mathbb{E}f(X)$ for concave f .

, worth taking a digression and talk about the conditions required for Jensen inequality to hold.

In case of f taking vector-valued inputs, f is convex \Rightarrow (if f has a Hessian) H is positive semi-definite. Moreover, if f is strictly convex ($H > 0$), then $Ef(x) = f(Ex)$ holds true $\Leftrightarrow X$ is constant.

We now have $Eg(\cdot)$ giving a lower bound on $l(\theta)$, to make the bound tight for a particular value

of θ , we demand the inequality becomes equality, thus Jensen inequality becomes equality as

we satisfy the condition $P(x = Ex) = 1$, namely, $\frac{P(x^{(i)}, z^{(i)})\theta}{Q_i(z^{(i)})} = \text{Const}$. Since $\sum_z Q_i(z^{(i)}) = 1$,

$$Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)}; \theta)}{\sum_z P(x^{(i)}, z; \theta)} = \frac{P(x^{(i)}, z^{(i)}; \theta)}{P(x^{(i)}; \theta)} = P(z^{(i)} | x^{(i)}; \theta)$$

thus, simply set Q_i 's to be the posterior distribution of $z^{(i)}$. With this tight lower-bound, we then maximize $Eg(\cdot)$ RHS w.r.t the parameters, the EM algorithm can then be written as:

Repeat until convergence?

(E-step) For each i , set $Q_i(z^{(i)}) := P(z^{(i)} | x^{(i)}; \theta)$.

(M-step). Set $\theta := \arg \max \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$

With the above described algorithm, we obtain $l(\theta^t) = \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})}$ in the t^{th} iteration, parameters $\theta^{(t+1)}$ are then obtained by explicitly maximizing RHS, thus:

$$l(\theta^{(t+1)}) \geq \underset{\substack{\downarrow \text{Jensen} \\ \text{[end of } t^{\text{th}} \text{ iteration]}}}{\sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})}} \geq \underset{\substack{\max \text{ RHS} \\ \theta}}{\sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})}} = l(\theta^{(t)}).$$

Hence, we showed that EM causes the likelihood to converge monotonically. One reasonable convergence test would be to check if the increase in $l(\theta)$ between successive iterations is smaller than some preset tolerance parameter.

If we define $J(Q, \theta) = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$, then EM can also be viewed as coordinate ascent on J , in which E-step maximizes J w.r.t Q and M-step maximizes J w.r.t θ .

• Derivation of EM applying to Gaussians

Let's go back and revisit our example of fitting a mixture of Gaussians. We now see it as an application to demonstrate the derivation of M-step. The E-step is easy by following the established framework and calculate $w_j^{(i)} = Q_i(z^{(i)}=j) = P(z^{(i)}=j | x^{(i)}; \phi, \mu, \Sigma)$. Here, $Q_i(z^{(i)}=j)$ denotes the probability of $z^{(i)}$ taking value of j under distribution Q , and now we view $\{w_j^{(i)}\}_{j=1}^{m_K}$ as constants.

Next, we seek to maximize the quantity $\sum_{i=1}^m \sum_{j=1}^k \alpha_i(z^{(i)}) \log p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)$ w.r.t ϕ, μ, Σ .

In our particular case setting, the quantity becomes $\sum_{i=1}^m \sum_{j=1}^k \alpha_i(z^{(i)}=j) \log p(x^{(i)}|z^{(i)}=j) p(z^{(i)}=j) / \alpha_i(z^{(i)}=j)$
 $= \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \left(\frac{1}{\sum_j} \exp(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)) \right) / w_j^{(i)}$. If we take derivative w.r.t. μ_j , we find
 $\nabla_{\mu_j} \mathcal{L}(\phi, \mu, \Sigma) = -\frac{1}{2} \nabla_{\mu_j} \sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) = \frac{1}{2} \sum_{i=1}^m w_j^{(i)} (\Sigma_j^{-1} x^{(i)} - \Sigma_j^{-1} \mu_j) := 0$

Solving this for μ_j therefore yields the update rule: $\mu_j := \frac{\sum_i w_j^{(i)} x^{(i)}}{\sum_i w_j^{(i)}}$.

If we take derivative w.r.t. ϕ_j , to deal with constraint that $\sum_j \phi_j = 1$, we construct the Lagrangian:

$$\mathcal{L}(\phi) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j + \beta \left(\sum_{j=1}^k \phi_j - 1 \right), \text{ where } \beta \text{ is the Lagrange multiplier.}$$

Notice we didn't explicitly constrain $\phi_j \geq 0$, cause this constraint is automatically satisfied by the solution

we are about to find. Moving on, $\frac{\partial}{\partial \phi_j} \mathcal{L}(\phi) = \sum_{i=1}^m \frac{w_i^{(i)}}{\phi_j} + \beta := 0$, i.e., $\phi_j \propto \sum_{i=1}^m w_i^{(i)}$.

Plugging in the constraint, we obtain $\sum_{j=1}^k \frac{m}{\sum_{i=1}^m w_i^{(i)}} = 1 \Rightarrow \beta = -\frac{m}{\sum_{j=1}^k \sum_{i=1}^m w_i^{(i)}} = -m$, thus the M-step update for parameters ϕ_j are $\phi_j := \frac{1}{m} \sum_{i=1}^m w_i^{(i)}$.

The derivation for update rule of Σ_j is also straight-forward, if we take the quantity and take derivative w.r.t. Σ_j , we obtain:

$$\begin{aligned} \nabla_{\Sigma_j} \sum_{i=1}^m w_i^{(i)} \log |\Sigma_j|^{-\frac{1}{2}} \exp(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)) &= \sum_{i=1}^m w_i^{(i)} \left[-\frac{1}{2} \Sigma_j^{-1} - \frac{1}{2} \nabla_{\Sigma_j} (\Sigma_j^{-1})^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right] \\ &= -\frac{1}{2} \sum_{i=1}^m w_i^{(i)} \left[\Sigma_j^{-1} - ((x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T) \Sigma_j^{-1} \cdot \Sigma_j^{-1} \right] \end{aligned}$$

By setting the derivative to zero and solving for Σ_j we yield the update rule $\Sigma_j := \frac{\sum_{i=1}^m w_i^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_i^{(i)}}$.

Factor Analysis

In our previous mixture of Gaussians setting, we imagine problems, where we have sufficient data to be able to discern the multiple-Gaussian structure in the data. For instance this would be the case when we have $m \gg d$ samples $m \gg$ dimension of features d .

Now, consider the setting in which $d \gg m$, how the m points span only a low-dimensional subspace of \mathbb{R}^d , if we try to model the data with a single Gaussian, then with MLE we would find

$$\Sigma_{MLE} = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{MLE})(x^{(i)} - \mu_{MLE})^T, \text{ which poses the fact that } \text{rank}(\Sigma) \leq \sum_{i=1}^m \text{rank}[(x^{(i)} - \mu)(x^{(i)} - \mu)^T]$$

$\leq m \ll d$, meaning that Σ_{MLE} is singular, which poses numerical problems when using multivariate Gaussian distribution, as a singular covariance matrix suggests that Gaussian places all of its probability in the affine space spaned by data.

To deal with the problem mentioned above, we posit a joint distribution on (x, z) as follows:

$$z \sim \mathcal{N}(0, I_{k \times k}), \quad x|z \sim \mathcal{N}(\mu + \Lambda z, \Psi). \quad \text{Parameters: } \mu \in \mathbb{R}^d, \Lambda \in \mathbb{R}^{d \times k}, \text{ diagonal matrix } \Psi \in \mathbb{R}^{d \times d}$$

with \mathbb{R}^k being latent variables. We can think (and thus choose) k to be the # "driving forces", namely, choose $k (< d)$ to be the number of things that might affect feature x .

The factor analysis model can be interpreted as: each datapoint $x^{(i)}$ is generated by sampling a k -dimension multivariate Gaussian $z^{(i)}$. Then it is mapped to a d -dimensional affine space of \mathbb{R}^d by performing linear operations $\mu + \Lambda z^{(i)}$ on $z^{(i)}$. Lastly, we add gaussian noise $\varepsilon \sim \mathcal{N}(0, \Psi)$ to it to obtain $\underline{x^{(i)}} = \mu + \Lambda z^{(i)} + \varepsilon$.

It is vital to recall that the basic thought for our model is $p(x, z) = p(x|z)p(z)$ with z being latent variables, so our random variables x, z have a joint Gaussian distribution $\begin{bmatrix} z \\ x \end{bmatrix} \sim \mathcal{N}(\mu_{zx}, \Sigma)$.

We observe that $\mathbb{E}[x] = \mathbb{E}[\mu + \Lambda z + \varepsilon] = \mu$, in order to find Σ , we should first note that

$$\text{Cov}\left(\begin{bmatrix} z \\ x \end{bmatrix}\right) = \Sigma \triangleq \begin{bmatrix} \Sigma_{zz} & \Sigma_{zx} \\ \Sigma_{xz} & \Sigma_{xx} \end{bmatrix} = \mathbb{E}\left[\begin{pmatrix} z - \mu_z & z - \mu_z \\ x - \mu_x & x - \mu_x \end{pmatrix} \begin{pmatrix} z - \mu_z & z - \mu_z \\ x - \mu_x & x - \mu_x \end{pmatrix}^T\right] = \mathbb{E}\left[\begin{pmatrix} (z - \mu_z)(z - \mu_z)^T & (z - \mu_z)(x - \mu_x)^T \\ (x - \mu_x)(z - \mu_z)^T & (x - \mu_x)(x - \mu_x)^T \end{pmatrix}\right]$$

$$\text{and our model gives } \Sigma_{zz} = \mathbb{E}[(z - \mu_z)(z - \mu_z)^T] = \text{Cov}(z) = I_{k \times k}, \quad \text{z, } \varepsilon \text{ is independent}$$

$$\Sigma_{zx} = \mathbb{E}[(z - \mu_z)(x - \mu_x)^T] = \mathbb{E}[z(x^T \mu^T)] = \mathbb{E}[zz^T] \Lambda^T + \mathbb{E}[z\varepsilon^T] = \Lambda^T$$

$$\Sigma_{xx} = \mathbb{E}[(x - \mu_x)(x - \mu_x)^T] = \mathbb{E}[\Lambda zz^T \Lambda^T + \Lambda z\varepsilon^T + \underbrace{\varepsilon\varepsilon^T \Lambda^T}_{\text{expectation zero}} + \varepsilon\varepsilon^T] = \Lambda \Lambda^T + \Psi.$$

We now see that the marginal distribution of x is given by $x \sim \mathcal{N}(\mu, \Lambda \Lambda^T + \Psi)$. Thus, given a training set

$$\{(x^{(i)})\}_{i=1}^m, \text{ the log likelihood of parameters is } \ell(\mu, \Lambda, \Psi) = \sum_{i=1}^m \log \frac{1}{(2\pi)^{d/2}} |\Lambda \Lambda^T + \Psi|^{-1/2} \exp\left(-\frac{1}{2} (x^{(i)} - \mu)^T (\Lambda \Lambda^T + \Psi)^{-1} (x^{(i)} - \mu)\right).$$

Maximizing this quantity wrt. parameters explicitly is hard, and we turn to EM algorithm for help.

• EM for Factor Analysis.

First, let's remind ourselves that the conditional of a joint Gaussian is Gaussian, suppose that

$$\overset{m}{\underset{n \uparrow}{\begin{bmatrix} x_A \\ x_B \end{bmatrix}}} \sim \mathcal{N}\left(\begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}\right), \text{ then the conditional densities are also Gaussian:}$$

$$x_A|x_B \sim \mathcal{N}(\mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (x_B - \mu_B), \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}), \quad x_B|x_A \sim \mathcal{N}(\mu_B + \Sigma_{BA} \Sigma_{AA}^{-1} (x_A - \mu_A), \Sigma_{BB} - \Sigma_{BA} \Sigma_{AA}^{-1} \Sigma_{AB}).$$

Using this property we can conclude $z^{(i)}|x^{(i)}, \mu, \Lambda, \Psi \sim \mathcal{N}(\mu_{z|x^{(i)}}, \Sigma_{z|x^{(i)}})$, with

$$\mu_{z|x^{(i)}} = \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} (x^{(i)} - \mu), \quad \Sigma_{z|x^{(i)}} = I - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \Lambda.$$

Plug these into $Q_i(z^{(i)})$ for EM we obtain (E-step)

$$Q_i(z^{(i)}) = (2\pi)^{-d/2} \left| \sum_{z^{(i)}|x^{(i)}} \right|^{-1/2} \exp\left(-\frac{1}{2} (z^{(i)} - \mu_{z|x^{(i)}})^T \Sigma_{z|x^{(i)}}^{-1} (z^{(i)} - \mu_{z|x^{(i)}})\right).$$

$$\text{In M-step we need to maximize } \sum_{i=1}^m \int_{\mathbb{R}^{d_z}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}|z^{(i)}; \mu, \Delta, \Psi)}{Q_i(z^{(i)})} dz^{(i)} = \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim Q_i} [\log p(x^{(i)}|z^{(i)}) - \log Q_i(z^{(i)}) + \log p(z^{(i)})]$$

Dropping out $Q_i(z^{(i)})$ which is determined in the E-step and $p(z^{(i)})$ which contains no parameter,

$$\text{we find that we need to maximize } \sum_{i=1}^m \mathbb{E}_{z^{(i)}} [\log p(x^{(i)}|z^{(i)}; \mu, \Delta, \Psi)]$$

$$= \sum_{i=1}^m \mathbb{E}_{z^{(i)}} \left[-\frac{d}{2} \log \pi - \frac{1}{2} \log |\Psi| - \frac{1}{2} (x^{(i)} - \mu - \Delta z^{(i)})^\top \Psi^{-1} (x^{(i)} - \mu - \Delta z^{(i)}) \right] \quad (1)$$

Maximize Eq(1) w.r.t Δ , we obtain:

$$\begin{aligned} & \nabla_{\Delta} \sum_{i=1}^m -\frac{1}{2} \mathbb{E}_{z^{(i)}} [(x^{(i)} - \mu - \Delta z^{(i)})^\top \Psi^{-1} (x^{(i)} - \mu - \Delta z^{(i)})] \\ &= \nabla_{\Delta} -\frac{1}{2} \sum_{i=1}^m \mathbb{E}_{z^{(i)}} [z^{(i)\top} \Delta^\top \Psi^{-1} \Delta z^{(i)} + (\mu - x^{(i)})^\top \Psi^{-1} \Delta z^{(i)} + (\Delta z^{(i)})^\top \Psi^{-1} (\mu - x^{(i)})] \\ &= \nabla_{\Delta} \sum_{i=1}^m \mathbb{E}_{z^{(i)}} \left[-\frac{1}{2} \text{tr}(\Delta^\top \Psi^{-1} \Delta z^{(i)\top} z^{(i)}) - \text{tr}(\Delta^\top \Psi^{-1} (\mu - x^{(i)}) z^{(i)\top}) \right] \\ &= \sum_{i=1}^m \mathbb{E}_{z^{(i)}} [-\Psi^{-1} \Delta z^{(i)\top} z^{(i)} + \Psi^{-1} (x^{(i)} - \mu) z^{(i)\top}] \stackrel{\text{set } 0}{=} 0 \Rightarrow \Delta = \left(\sum_{i=1}^m (x^{(i)} - \mu) \mathbb{E}_{z^{(i)}} [z^{(i)\top} z^{(i)}] \right)^{-1} \end{aligned}$$

And with our definition Q_i being Gaussian with mean $\mu_{z^{(i)}|x^{(i)}}$ and covariance $\Sigma_{z^{(i)}|x^{(i)}}$, we find that

$$\mathbb{E}_{z^{(i)} \sim Q_i} [z^{(i)\top}] = \mu_{z^{(i)}|x^{(i)}}, \quad \mathbb{E}_{z^{(i)} \sim Q_i} [z^{(i)\top} z^{(i)}] = \mu_{z^{(i)}|x^{(i)}} \mu_{z^{(i)}|x^{(i)}}^\top + \Sigma_{z^{(i)}|x^{(i)}}$$

Similarly by taking derivative w.r.t. Ψ and μ we obtain:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}, \quad \Psi = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)\top} - \frac{1}{m} \sum_{i=1}^m \mu_{z^{(i)}|x^{(i)}} \mu_{z^{(i)}|x^{(i)}}^\top - \Delta \mu_{z^{(i)}|x^{(i)}} x^{(i)\top} + \Delta (\mu_{z^{(i)}|x^{(i)}} \mu_{z^{(i)}|x^{(i)}}^\top + \Sigma_{z^{(i)}|x^{(i)}}) \Delta^\top$$

Notice μ doesn't change as the parameters are varied, this can be calculated just once, and Ψ is a diagonal matrix as we demand, so Ψ can be found by setting $\Psi_{ii} = \Phi_{ii}$ ($i=1, 2, \dots, d$) (i.e. Ψ constraints only the diagonal entries of Φ). (Q: Would this selection affect the optimal solution? or can proof Φ is already diagonal if ignoring precision issues?)

Principal Components Analysis

Like Factor Analysis, PCA view high-dimensional data as lying on low dimension subspace with some noise. Prior to running PCA per se, typically we first preprocess the data $x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j}$ where $\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$ and $\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n (x_j^{(i)} - \mu_j)^2$. This normalization rescales the different attributes to make them more comparable and may be omitted if we had a priori knowledge that different attributes are all on the same scale. (e.g. grayscale image).

We propose two ways of deriving the idea of PCA. First, notice our basic view of a low dimension subspace, we want to choose a subspace that is "closest" to all the data points; Another perspective is to choose a unit vector u , and when we project the data onto the direction corresponding to u ,

as much as possible of the variance in the data is retained. These two interpretations can be related as projection length and projection height is related by $\|x^{(i)}\|_2^2$, and lead to the optimization problem (here we use the second intuition and maximize the length of projection, which, given unit vector u and point x , is given by $\frac{\langle x, u \rangle}{\|u\|} = x^T u$): $\max_u \frac{1}{n} \sum_{i=1}^n (x^{(i)T} u)^2$ s.t. $\|u\|_2 = 1$.

To tackle with the optimization problem we construct the Lagrangian $\mathcal{L}(u) = \frac{1}{n} \sum_{i=1}^n (x^{(i)T} u)^2 + \beta(\|u\|_2 - 1)$ $= \frac{1}{n} \sum_{i=1}^n u^T x^{(i)} x^{(i)T} u + \beta(u^T u - 1) = u^T \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)T} \right) u + \beta(u^T u - 1) = u^T \Sigma u + \beta(u^T u - 1)$, with Σ being the empirical covariance matrix of the data. (Recall that normalization step ensures zero mean). It is easy to find $\nabla_u \mathcal{L}(u)$ and by setting this to zero we see that u should be chosen as eigenvectors of Σ^* . If we wish to project our data into a k -dimensional subspace ($k < d$), we should choose u_1, u_2, \dots, u_k to be the top k eigenvectors of Σ , and since Σ is symmetric, u_1, u_2, \dots, u_k can be transformed to be orthogonal by comparison of corresponding eigenvalue.

basis of k -dimension subspace in the original d -dimensional sample space. The original data can now be modelled as $\{x^{(i)}\}_{i=1}^n \rightarrow (u_1^T x^{(i)}, u_2^T x^{(i)}, \dots, u_k^T x^{(i)}) \stackrel{?}{=} y^{(i)} \in \mathbb{R}^k$. [dimensionality reduction from $\mathbb{R}^d \rightarrow \mathbb{R}^k$.]

If we wish to "uncompress" the data and model $x^{(i)}$ from $y^{(i)}$, we should be aware of the fact that $x^{(i)} \approx y_1^{(i)} u_1 + y_2^{(i)} u_2 + \dots + y_k^{(i)} u_k$, with approximation coming from the omission of noises.

Application Analysis

- Scenarios considered to use PCA:

1. Visualization. By reducing data to 2/3 dimension, we can plot the $y^{(i)}$'s to visualize the data, during which process we might be able to observe interesting features such as clusters (indicating types of the features) or trajectory of how features evolve.
2. Compression for Efficiency. It turns out most of high-dimensional data lies on low-dimensional subspace, which is the assumption that led us derive Factor Analysis & PCA (and also a fact of life).

Apart from improving computational efficiency, reducing data's dimension can also reduce the

complexity of hypothesis class considered (lower dimensional input probably have lower VC dimension).

- Questionable scenarios and ideas:

1. Reduce Overfitting. This idea comes from lower feature numbers (leads to lower variance, but this technique works as often as it fails, consider using regularization instead).

2. Outlier Detection. In early stages of face recognition, reducing original image data $x^{(i)} \in \mathbb{R}^{100 \times 100}$ to lower dimensions (say 2) and plotting the processed data results in eigenfaces, and measuring $\|y^{(i)}_1 y^{(j)}_2\|_2^2$ results in a surprisingly good face-matching algorithm. In eigenfaces, it seems the principal components can be interpreted (e.g. glasses, mustache) and PCA helps reduce noise on feature vectors, but **individual eigenvector is quite noisy so avoid examining single eigenvector**, in which case interpretation lacks theoretical proof and is a product of researchers' imagination.

- General Rule of Thumb: For implementations using PCA, consider how the algorithm does on original data to avoid adding PCA unnecessarily.

[positions (numerical) of individual eigenvector is unstable to perturbations in data]
[But subspace spanned by these eigenvectors is quite stable]

A brief conclusion of our tour in unsupervised learning algorithms:

| | | |
|-----------|--|--|
| | Model $p(x)$ (e.g. For anomaly detection) | Non-probabilistic (e.g. compression, visualization) |
| Subspaces | Factor Analysis | PCA |
| Clusters | Mixture of Gaussians (not necessarily clusters) | k-means |

(latent variables)

Apply EM since EM models $p(x)$ given $p(x_i)$.

Independent Components Analysis

The motivation for ICA is to extract independent sounds from their combinations, say $S \in \mathbb{R}^d$ with $s_j^{(i)}$ being j^{th} speaker's voice at time i ; what we observe is $x^{(i)} = As^{(i)}$ with $x_j^{(i)}$ being observation at j^{th} microphone, (for now we assume # speakers = # microphone) and A an unknown matrix $\in \mathbb{R}^{d \times d}$. Repeated observations gives a dataset $\{x^{(i)}\}_{i=1}^n$ and our goal is to recover the sources $s^{(i)}$ that had generated our data $x^{(i)}$. We will denote $W = \begin{bmatrix} -w_1^T \\ \vdots \\ -w_d^T \end{bmatrix} = A^{-1}$, then j^{th} source can be recovered as $s_j^{(i)} = w_j^T x^{(i)}$.

• ICA Ambiguities

Under what circumstances can $W = A^{-1}$ be recovered? There are some inherent ambiguities in A that are impossible to recover, given only the $x^{(i)}$'s.

First, it's trivial to see that there is no way to recover the correct scaling of the w_j 's. If a

single column of A were scaled by α , and the corresponding source was scaled by $\sqrt{\alpha}$, then there is no way to determine that this has happened given only the $x^{(i)}$'s. Notice that the scaling factor controls the volume of that speaker's voice, this ambiguity would not matter since we are only trying to identify the independent speakers. (A is the case for most ICA applications).

Also notice given only the $x^{(i)}$'s, it is impossible to distinguish between TW and PW , with P being any permutation matrix $\in \mathbb{R}^{d \times d}$. Intuitively, there are no "absolute" labeling for the speakers.

It is vital to point out the ambiguity with sources being generated from Gaussian. For illustration, let's say $s_i \sim \mathcal{N}(0, I_{d \times d})$, so the contours are circles ($d=2$ case) centered on the origin. Now, suppose we observe some $x = As \sim \mathcal{N}(0, AA^T)$, and let R be an arbitrary orthogonal matrix, and let $A' = AR$. The distribution of $x' = A's$ is still $\mathcal{N}(0, AA^T)$. Thus, there is no way to tell if sources were mixed using A or A' , this is because the density of Gaussian is rotationally symmetric*. However, for natural voices, the data turns out not being Gaussian, with high/low pitches frequently observed but would be treated as outliers in Gaussian.

• Algorithm Derivation

We suppose that the distribution of each source s_j is given by density p_s , now given s_j are iid we obtain $p(x) = \prod_{j=1}^q p_s(s_j) \Rightarrow \prod_{j=1}^q p_s(w_j^T x) | A^{-1}| = p(x)$. [recall that $x = As = w^T s$ gives $p(x) = p_s(wx) |Tw|$].

A "default" choice of density p is $p_s(s) = \sigma'(s)$ with σ being sigmoid function, which, when chosen as CDF, gives a PDF that has "fatter" tails than the Gaussian.

Note: Default choice implies that we had no prior knowledge of the sources' densities; if indeed we have, then that density would be a good substitute of default. The Laplace distribution $p(x) = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right)$ is also a reasonable choice.

The representation here assumes that either the data $x^{(i)}$ has been preprocessed to have zero mean, or that it can naturally be expected to have zero mean (such as acoustic signals), cause under the condition $\sigma'(s)$ is symmetric, $E[s] = 0 \Rightarrow E[x] = E[As] = 0$.

Moving on to MLE, the log likelihood is given by

$$P(TW) = \sum_{i=1}^n \left[\sum_{j=1}^q \log p_s(w_j^T x^{(i)}) + \log |Tw| \right] = \sum_{i=1}^n \left[\sum_{j=1}^q \log \sigma'(w_j^T x^{(i)}) + \log |Tw| \right].$$

Taking derivative of $P(TW)$ gives us the stochastic gradient ascent update:

$$W := W + \alpha \nabla_W f(W) = W + \alpha \left(\begin{bmatrix} 1 - \alpha(w^T x^{(i)}) \\ \vdots \\ 1 - \alpha(w^T x^{(i)}) \end{bmatrix} x^{(i)T} + (w^{-1})^T \right)$$

where α is the learning rate. After the algorithm converges, we then compute $S^{(i)} = W x^{(i)}$ to recover sources. When deriving likelihood of data, we implicitly assumed that $x^{(i)}$'s were independent of each other, which is clearly incorrect for speech and other time series data. But it can be shown that having correlated training examples will not hurt performance of algorithm if we have sufficient data. However, for problems where successive training examples are correlated, it sometimes helps SGD to converge if we run stochastic gradient descent on a randomly shuffled copy of training set.

A brief mentioning of #speakers \neq # microphones scenarios, first note if #speakers $<$ # microphones, then we can assume there are silent speakers, but #speakers $>$ # microphones is still a cutting edge research area, with some findings suggest it's possible to distinguish speakers if pitches of voice differ to some extent.

Boosting

Roughly, the idea of boosting is to take a weak learning algorithm (one that is slightly better than random guessing) and transform it into a strong classifier. Boosting procedures proceed by taking a collection of weak classifiers, and then reweighting their contributions to form a classifier with much better accuracy than any individual classifier.

Let's formulate the problem. Assume we have raw inputs $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$, $x^{(i)} \in \mathbb{R}^n$, $y^{(i)} = \{-1, +1\}$. as usual in a binary classification setting. We also assume we have an infinite collection of feature functions $\phi_j: \mathbb{R}^n \rightarrow \{-1, 1\}$ and an infinite weight vector $\theta = (\theta_1, \theta_2, \dots)^T$ with finite non-zero entries. Our hypothesis $h_\theta(x) = \text{sign}(\sum_{j=1}^m \theta_j \phi_j(x)) \triangleq \text{sign}(\theta^T \phi(x))$. During process of boosting, we start by assigning equal weight to each training example and continuously perform reweighting where examples that it misclassifies receive higher weight and examples it correctly classified receive lower weight. We now formalize the concept of "weight" to be a distribution on the examples where $\sum_{i=1}^m p^{(i)} = 1$ and $\forall i, p^{(i)} \geq 0$. We say that there is a weak learner with margin $\gamma > 0$ if for any distribution $p = (p^{(1)}, p^{(2)}, \dots, p^{(m)})$ on the m training examples there exists some weak hypothesis ϕ_j st.

$$\sum_{i=1}^m p^{(i)} \mathbb{1}\{y^{(i)} \neq \phi_j(x^{(i)})\} \leq \frac{1}{2} - \gamma. \quad (1)$$

That is, we assume there is some classifier that does slightly better than random guessing.

Now let's focus on deriving a loss function, notice:

$$\mathbb{1}\{\text{sign}(\theta^T \phi(x)) \neq y\} = \mathbb{1}\{y \theta^T \phi(x) \leq 0\} \leq \exp(-y \theta^T \phi(x)).$$

Thus our misclassification error $\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{\text{sign}(\theta^T \phi(x^{(i)})) \neq y^{(i)}\} \leq \frac{1}{m} \sum_{i=1}^m \exp(-y^{(i)} \theta^T \phi(x^{(i)}))$, which leads to assigning the loss function as $J(\theta) = \frac{1}{m} \sum_{i=1}^m \exp(-y^{(i)} \theta^T \phi(x^{(i)}))$. Notice the convenience of "arg min" operation on exponential loss, we perform coordinate descent update for risk $J(\theta)$. Suppose we wish to update coordinate k , define $w^{(i)} = \exp(-y^{(i)} \sum_{j \neq k} \theta_j \phi_j(x^{(i)}))$, and $W^+ := \sum_{i: y^{(i)} \phi_k(x^{(i)}) = 1} w^{(i)}$, $W^- := \sum_{i: y^{(i)} \phi_k(x^{(i)}) = -1} w^{(i)}$. Thus:

$$\arg \min_{\theta_k} J(\theta) = \arg \min_{\theta_k} \{ W^+ e^{\theta_k} + W^- e^{-\theta_k} \} = \frac{1}{2} \log \frac{W^+}{W^-}$$

Algorithm Derivation

We assume that we have access to a weak-learning algorithm, which at iteration t takes as an input a distribution p on the training set and returns a hypothesis ϕ_t satisfying margin condition (1), then the set of weak hypotheses returned by the weak learning algorithm can be denoted as $\{\phi_1, \phi_2, \dots, \phi_T\}$, and the full boosting algorithm can be described as:

For each iteration $t=1, 2, \dots$:

(1) Define weights $w^{(i)} := \exp(-y^{(i)} \sum_{t=1}^{t-1} \theta_t \phi_t(x^{(i)}))$. // current weak hypothesis: $\{\phi_1, \phi_2, \dots, \phi_{t-1}\}$.

and distribution $p^{(i)} := w^{(i)} / \sum_{j=1}^m w^{(j)}$. // idea of assigning bigger weights to misclassified examples

(2) Construct a weak hypothesis $\phi_t: \mathbb{R}^n \rightarrow \{-1, +1\}$ from the distribution

$p = (p^{(1)}, p^{(2)}, \dots, p^{(m)})$ on the training set [output from (1)]

(3) Compute $W_T^+ := \sum_{i: y^{(i)} \phi_t(x^{(i)}) = 1} w^{(i)}$ and $W_T^- := \sum_{i: y^{(i)} \phi_t(x^{(i)}) = -1} w^{(i)}$

and set $\theta_t := \frac{1}{2} \log \frac{W_T^+}{W_T^-}$.

Recall we construct our classifier as $h_0(x) = \text{sign}(\sum_{j=1}^m \theta_j \phi_j(x))$ and said θ only have finite non-zero entries.

We now argue that the boosting procedure * achieves 0 training error * and provide a * rate of convergence to zero *.

Convergence of Boosting

Lemma. Let

$$\hat{J}(\theta) = \frac{1}{m} \sum_{i=1}^m \exp(-y^{(i)} \sum_{t=1}^T \theta_t \phi_t(x^{(i)}))$$

Then $\hat{J}(\theta^{(t)}) \leq \sqrt{1-4r^2} \hat{J}(\theta^{(t-1)})$, with r being margin specified by condition (i).

Proof. We inspect that $\hat{J}(\theta^{(t)}) = \min \{W_t^+ e^{\theta}, W_t^- e^{-\theta}\} = 2\sqrt{W_t^+ W_t^-}$.

$$\text{While } \hat{J}(\theta^{(t-1)}) = \frac{1}{m} \sum_{i=1}^m \exp(-y^{(i)} \sum_{t=1}^{T-1} \theta_t \phi_t(x^{(i)})) = W_t^+ + W_t^-.$$

$$\text{We know from condition (i) that } \sum_{i=1}^m p^{(i)} \mathbb{1}\{y^{(i)} \neq \phi_t(x^{(i)})\} \leq \frac{1}{2} - r \xrightarrow[\text{of } p^{(i)}]{\text{algorithm specification}} \frac{W_t^+}{W_t^+ + W_t^-} \leq \frac{1}{2} - r.$$

$$\text{Rewriting this expression we get } W_t^- \leq (\frac{1}{2} - r)(W_t^+ + W_t^-) \text{ or } W_t^+ \geq \frac{1+2r}{1-2r} W_t^-.$$

$$\text{which indicates } \lambda \text{ (in the minimum defining } \hat{J}(\theta^{(t)})) = \frac{1}{2} \log \frac{W_t^+}{W_t^-} \geq \frac{1}{2} \log \frac{1+2r}{1-2r}.$$

$$\begin{aligned} \text{Leading to } \hat{J}(\theta^{(t)}) &\leq W_t^+ \sqrt{\frac{1+2r}{1-2r}} + W_t^- \sqrt{\frac{1-2r}{1+2r}} \\ &= W_t^+ \sqrt{\frac{1+2r}{1+2r}} + (1-2r) W_t^- \sqrt{\frac{1+2r}{1+2r}} + 2r \sqrt{\frac{1+2r}{1+2r}} W_t^- \\ &\leq W_t^+ \left(\sqrt{\frac{1+2r}{1+2r}} + 2 \sqrt{\frac{1-2r}{1+2r}} \right) + (1-2r) \sqrt{\frac{1+2r}{1-2r}} W_t^- \\ &= (1-4r^2)(W_t^+ + W_t^-) = (1-4r^2) \hat{J}(\theta^{(t-1)}). \end{aligned}$$

□

We initialize the training procedure at $\theta^{(0)} = \vec{0}$ making the initial empirical risk 1. Notice

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{\text{sign}(\theta^T \phi(x^{(i)})) \neq y^{(i)}\} \leq \hat{J}(\theta)$$

Thus if $\hat{J}(\theta) < \frac{1}{m}$ then θ makes no mistakes on the training data*. After T iterations of boosting the empirical risk satisfies $\hat{J}(\theta^{(T)}) \leq (1-4r^2)^{\frac{T}{2}} \hat{J}(\theta^{(0)}) = (1-4r^2)^{\frac{T}{2}}$. We demand $\hat{J}(\theta^{(T)}) < \frac{1}{m}$ to find the number of iteration needed: $T > \frac{-\log m}{-\log(1-4r^2)}$. Notice RHS $\leq \frac{\log m}{2r^2}$, thus with more than $\frac{\log m}{2r^2}$ rounds of boosting we can come up with a classifier of zero training loss. Although it is also important for the model to have low generalization error, we will ignore test part and focus on training error in this chapter.

Implementing Weak-learners

There are a number of strategies for weak learners and in this section we focus on one, known as decision stumps. For concreteness of discussion, let's suppose input variables $x \in \mathbb{R}^n$ are real-valued.

A decision stump is parameterized by a threshold s and index j , namely select j^{th} entry of x and compare x_j with s , and returns $\phi_j(x) = \text{sign}(x_j - s) = \begin{cases} +1, & \text{if } x_j \geq s \\ -1, & \text{otherwise.} \end{cases}$

Recall our goal here is, given a distribution $(p^{(1)}, p^{(2)}, \dots, p^{(m)})$ on the training set, we wish to choose a decision stump of the above form to minimize error on training set, that is, find a threshold $s \in \mathbb{R}$ and feature index j s.t.

$$\widehat{\text{Err}}(\phi_{j,s}, P) = \sum_{i=1}^m P^{(i)} \mathbb{1}\{\phi_{j,s}(x^{(i)}) \neq y^{(i)}\} = \sum_{i=1}^m P^{(i)} \mathbb{1}\{y^{(i)} \text{sign}(x_j^{(i)} - s) = -1\} = \sum_{i=1}^m P^{(i)} \mathbb{1}\{y^{(i)}(x_j^{(i)} - s) \leq 0\}.$$

is minimized. Notice the only values s for which training error can change are the values $x_j^{(i)}$, if for each feature $j=1, 2, \dots, n$, we sort the raw input features s.t. $x_j^{(1)} \geq x_j^{(2)} \geq \dots \geq x_j^{(m)}$, then calculation of the above training error can be rewritten as $\sum_{k=1}^m P^{(ik)} \mathbb{1}\{y^{(ik)}(x_j^{(ik)} - s) \leq 0\}$. Choose for s , iterating over $\#$ training samples and $\#$ features let us choose a index j and threshold s that gives the best decision stump.

One important thing to note is that $\widehat{\text{Err}}(\phi_{j,s}, P) = 1 - \widehat{\text{Err}}(-\phi_{j,s}, P)$, thus it is important to also track the smallest value of $1 - \widehat{\text{Err}}(\phi_{j,s}, P)$ over all thresholds, because it may be smaller than $\widehat{\text{Err}}(\phi_{j,s}, P)$ which gives a better weak learner.

Now we discuss about variations on the basic boosted decision stumps. First, we do not need input features x_j be real-valued, let's say some x_j might be categorical with $x_j \in \{1, 2, \dots, k\}$, in which case decision stumps can take the form $\phi_j(x) = \begin{cases} 1, & \text{if } x_j = 1 \\ -1, & \text{otherwise} \end{cases}$ as well as variants setting $\phi_j(x) = 1$ if $x_j \in C$ for some set $C \subseteq \{1, 2, \dots, k\}$. Another natural variation is the boosted decision tree, in which instead of a single level decision-making, we consider conjunctions of trees of decisions.

Reinforcement Learning

Note for many sequential decision making and control problems, there is an absence of (or difficult to come up with) explicit supervision (e.g. labeling). In the reinforcement learning framework, we will instead provide our algorithm only a reward function, which would be an indicator for whether the algorithm is doing well, and it will be the algorithm's job to choose actions over time so as to obtain large rewards.

• Markov Decision Process

MDP formalizes our reinforcement learning setting by proposing a tuple $(S, A, \{P_{sa}\}, \gamma, R)$, where:

- * S : set of states (e.g. in a chess game setting, S will be all possible chess positions)

- * A : set of actions (e.g. the set of all possible moves we can adapt)

- * P_{sa} : state transition probabilities. For each state $s \in S$ and $a \in A$, P_{sa} is a distribution over the state space that we can transit to, for all states s' we can arrive at, $\sum_{s'} P_{sa}(s') = 1$.

- * γ : discount factor in $[0, 1]$.

* R : reward function mapping $S \times A$ to \mathbb{R} , if written as a function of state S only, $R: S \mapsto \mathbb{R}$.

The dynamics of an MDP proceeds as follows: Start in some state s_0 and get to choose a move $a_0 \in A$, as a result of our choice, the state of MDP randomly transitions to some successor state s_1 , drawn according to $P_{s_0 a_0}$, pictorially, $s_0 \xrightarrow{a_0} s_1 \sim P_{s_0 a_0} \xrightarrow{a_1} s_2 \sim P_{s_1 a_1} \xrightarrow{a_2} s_3 \sim P_{s_2 a_2} \dots$

Upon visiting the sequence of states s_0, s_1, \dots with actions a_0, a_1, \dots , our total payoff is given by

$$R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots$$

Let's interpret the existence of discount factor, intuitively, with γ^t monotonically decreasing, we are encouraging the algorithm to get the reward (positive) as quickly as possible, pragmatically, since we are defining payoff in terms of series (which is the case when there is no absorbing state), by adding γ^t , with Dirichlet test of convergence, our series will surely converge.

[One finicial interpretation is seeing γ as the value of money, which depreciate over time] ^(interest rate)

Our goal in reinforcement learning is to choose actions over time so as to maximize the expected value of the total payoff $\mathbb{E}[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$.

We now introduce policy as a function $\pi: S \mapsto A$, executing policy π indicates whenever we are in state s , we take action $a = \pi(s)$. Define the value function for a certain policy π :

$$V^\pi(s) = \mathbb{E}[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | s_0 = s, \pi] \quad \text{though not a random variable,}$$

which is simply the expected payoff upon starting in state s and act according to π . ^{we will note as condition on π .}

Note that $V^\pi(s) = \mathbb{E}[R(s_0) + \gamma V^\pi(s_1)]$ and $s_1 \sim P_{s_0 \pi(s_0)}$ which gives rise to Bellman equation:

$$V^\pi(s) = \underbrace{R(s_0)}_{\text{immediate reward}} + \gamma \sum_{s' \in S} \underbrace{P_{s \pi(s)}}_{\text{prob of transition}} \underbrace{V^\pi(s')}_{\text{expected sum of discounted rewards after first step}}$$

With Bellman equation, given a specific policy π in a finite-state MDP ($|S| < \infty$), we can write down one such equation for each $s \in S$. This gives us a set of $|S|$ linear equations in $|S|$ variables, which can be efficiently solved.

We naturally ask what is the ^{best} possible value of expected sum for any policy? Leading to defining

$$V^*(s) = \max_{\pi} V^\pi(s).$$

and a natural variation of Bellman's equations

$$V^*(s) = R(s) + \max_{a \in A} \gamma \sum_{s' \in S} P_{s a} V^*(s') \quad \begin{array}{l} \text{optimal value function} \\ \text{in the following process} \end{array} \quad (2)$$

↓
optimal action current step

We also define a policy $\pi^*: S \rightarrow A$ as follows:

$$\pi^*(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s') \quad (3)$$

which is the action that attains the maximum in the "max" in Eq(2), with extrapolation to S , π^* is the optimal policy. Note that π^* has the property of being optimal policy for all states $s \in S$, this means we can use the same policy π^* regardless of initial state of our MDP.

Solving finite-state MDPs

With the framework proposed in the previous section, our roadmap is now (1) find algorithm to compute $V^*(s) = V^{*\pi}(s)$, (2) use Eq(3) to find the optimal policy. For now, we will consider only MDPs with $|S| < \infty$ and $|A| < \infty$. In this section we will also assume that the state transition probabilities $\{P_{sa}\}$ and reward function R is known.

* Algorithm: Value Iteration

1. For each state s , initialize $V(s) := 0$ // then $[V(s_1), V(s_2), \dots, V(s_n)]^T \in \mathbb{R}^{|S|}$ with all entries being 0

2. Repeat until convergence. {

For every state s , update $V(s) := R(s) + \max_{a \in A} \gamma \sum_{s' \in S} P_{sa}(s') V(s')$ }

There are two possible ways of performing updates in the inner loop,

i) Synchronous update. Here RHS is always going to be old estimation value, and we are

computing RHS for $|S|$ times and then synchronously update all entries of $[V(s_1), V(s_2), \dots, V(s_n)]^T$.

We can proof, Value iteration causes $[V(s_1), V(s_2), \dots, V(s_n)]^T$ to converge to V^* in exponential speed.

ii) Asynchronous update. Here we would loop over the vector in some order and update the value

one at a time, so later computed $V(s_k)$ will be using / depending on refreshed value of $V(s_{k-1})$ (if

s_k is the successive state of s_{k-1}).

* Algorithm: Policy Iteration

1. Initialize π randomly. // for each of the states, pick a random action

2. Repeat until convergence. { // actually policy iteration * gets to exactly $V^{*\pi}$ (like normal equations in LR)

a) Solve the value function V^π using the set of Bellman's equations, and set $V := V^\pi$.

b) For each state s , let $\pi(s) := \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s')$ }, // pretend V current is optimal \Rightarrow greedy policy.

Notice step 2.(a) in policy iteration is done at the expense of inverting a matrix, so for small problems

(i.e. problems with small $|S|$)

policy iteration is a favorable choice, when in large problems, prioritize using value iteration.

After solving for V^* , if using policy iteration, then optimal policy is attained when algorithm converges, but if using value iteration, there is one more step to use Eq(3) and come up with π^* .

Practically, when we have absorbing states (e.g. winning/lossing, target) in our problem setting, we can model these states by *assigning their transition probabilities* to zero, therefore when using value/policy iteration, $V(\text{absorbing state}) = R(\text{absorbing state})$ which is the initial immediate reward, and won't alter as algorithm runs. For states that are not absorbing state, we would often *assign a small negative penalty* to each of those states (e.g. -0.02) to encourage algorithm to get the final reward as quickly as possible.

Learning a Model for MDP

So far we have discussed MDPs assuming that both state transition probability and reward function are known, but in realistic problems, these might not be prior knowledge and have to be estimated from empirical data.

[S, A, Y are usually known, and reward can be assigned]

In order to do so, let's say our MDP consists of a number of trials, we can then perform MLE to get state transition probabilities:

$$P_{sa}(s') = \frac{\text{(denominator)} \# \text{times we took action } a \text{ in state } s \text{ and got to } s'}{\# \text{times we took action } a \text{ in state } s}$$

Or, if numerator is 0 — corresponding to the case of never taking action a in state s , then we might simply estimate $P_{sa}(s')$ to be $Y_{1,1}$. It turns out we can apply Laplace smoothing here, like we did in naive Bayes, but we don't have to, since these solvers for MDP is not sensitive to zero values, unlike naive Bayes where having zero probability will be really problematic.

Using a similar procedure, if R is unknown, we can pick our estimate of the expected immediate reward $R(s)$ to be the average reward observed in state s .

Putting it together, we propose one possible algorithm for learning in an MDP with unknown state transition probabilities.

- * Algorithm:
1. Initialize π randomly // generate a random policy
 2. Take actions w.r.t. π for some number of trials to get experience in MDP.

Repeat {

(a) Execute π for certain time

(b) Using accumulated experience, update estimates for P_{sa}

[recall P_{sa} is estimated to be fraction of counts]

(c) Apply value/policy iteration with estimated P_{sa} to get a new estimated value function V .

(d). Update π to be the greedy policy w.r.t. V , namely, $\pi(s) = \arg \max_a \sum_{s'} P_{sa}(s')V(s')$.

One simple optimization is where we apply value iteration, instead of initializing with $V_i = \vec{0}$, we can initialize it with the solution found during the previous iteration of our algorithm. This simple technique will help value iteration converge more quickly by providing a much better initial starting point.

In some problem settings, the reward might be a random function of the state (e.g. stock price, self driving), and as mentioned before, this process can also help us estimate the immediate reward $R(s)$.

Exploration-Exploitation Tradeoff

The exploration vs. exploitation problem poses a question of how greedy should we be at just taking actions to maximize rewards, one simple setting for us to observe the problem would be:

| | | | | | |
|---|-----|-------|---|-----|-----|
| 5 | | ↑ | ↑ | ↑ | 0.9 |
| 4 | | -<0.1 | 0 | >+1 | |
| 3 | | ↓ | ↓ | ↓ | |
| 2 | | | | | |
| 1 | +10 | | | | |

exploration of known information

exploration for better reward = [+10 reward wasn't there for advance, need to figure out].

Say the initial position of our robot is (4,4) and the initial random policy leads our robot to find reward +1 at (6,4), and this would lead our algorithm to always go to (6,4), which, formally put, is locally greedy choices (lead to local optima instead of global optima). Also, when heading to reward at (6,4), our robot may wander off a bit and gets more information about right plane, and leave the left plane mostly unexplored.

Let's revise our algorithm, what we suggested is taking policies that is expected to maximize the total payoff, which is indeed locally greedy. What our algorithm did not do is exploration, which is, taking actions that seem less optimal currently (e.g. in our example, it would be going left multiple times just to see what happens) to gain more knowledge about the task.

It turns out the exploration vs. exploitation problem is not only academic, but occurs in practical settings. A typical example would be online ads company's strategy to pose ads. By posing ads the user is most likely to click on would drive short time interest, but might provide less knowledge about the user. If, say we are showing ads that, by our current knowledge of user interest, is not the ads that the user appears to be most interested in, but ads from other categories. By doing so, we might gain more

[chance for getting ad of interest or not?]

information about user interest and thus increase the effectiveness of these companies finding relevant ads.

One strategy to tackle this issue is a exploration strategy called ϵ -greedy, which instead of taking actions w.r.t. π , we now hold $(1-\epsilon)$ chance of taking actions $\xrightarrow{\text{greedy action}}$ w.r.t. π and ϵ chance of taking actions $\xrightarrow{\text{random action}}$ randomly. With this strategy, every now and then our robot might get to (1,1) and end up exploring the state space more thoroughly. It turns out, with ϵ -greedy exploration and sufficient amount of time, our algorithm will converge to the optimal policy (global optimum) for any discrete state MDP. A common heuristic for choosing ϵ is, start with a relatively large ϵ and gradually decrease it.

Another common exploration strategy is Boltzmann exploration, where we would choose probabilistically according to expected rewards, for example we would choose state s with probability $\exp(-P(s)/T)$ where $P(\cdot)$ is the payoff probability, and would decrease T (temperature) over time, with high T alluding exploration and low T leading to increasing bias toward best action.

Continuous State MDPs

For MDPs that have an infinite number of states, for example, when modeling a helicopter flying in 3d space, the states take the form of $(\overset{\text{roll}}{x}, \overset{\text{pitch}}{y}, \overset{\text{yaw}}{z}, \overset{3d \text{ coordinates}}{\phi}, \overset{\text{orientation}}{\theta}, \overset{\text{position}}{\psi}, \overset{\text{angle of pole}}{\dot{x}}, \overset{\text{angle of pole}}{\dot{y}}, \overset{\text{angle of pole}}{\dot{z}})$. Another example would be the inverted pendulum, where states take the form of $(x, \overset{\text{angle of pole}}{\theta}; \overset{\text{position}}{\dot{x}}, \overset{\text{position}}{\dot{\theta}})$, generally. In problems where $S = \mathbb{R}^d$, some of our former algorithm don't work well, or won't work at all.

In this section we will focus on solving continuous state MDPs.

Strategy 1. Discretization

This is perhaps the simplest thought, and usually works well for 1d and 2d problems (with the advantage of being simple and quick to implement) even without cleverly choosing the way to discretize. For example in 2d problem, we can use a grid to discretize the state space, and approximate the continuous-state MDP via a discrete-state one $(\bar{S}, \bar{A}, \{\bar{P}_{\bar{s}\bar{a}}\}, \bar{Y}, \bar{R})$.

There are two main downsides of this strategy. First, it naively assumes the value function takes a constant value over each discretization intervals. To better understand the limitation, consider fitting a piecewise constant function $\xrightarrow{\text{to}}$ a linear model, this "stair" function just isn't a good representation for many smooth functions. We might need really small intervals to get a good representation.

The second downside is called the curse of dimensionality. Suppose $S = \mathbb{R}^d$ and each dimension is

discretized into k values, then we have a total of k^d discrete states, and grows exponentially as problem dimension increases.

As a rule of thumb, discretization works well for 1d and 2d problems. With problems of 3d and 4d, we may have to decide which dimension affects the outcome the most, and discretize this influential state to more discrete states, and discretization might work in this scenario. But discretization very rarely works for problems with $\geq 6d$.

Strategy 2. Value Function Approximation

An alternative method would be to approximate V^* directly, like what we did in (linear) regression models, but here we will add some modifications and model $V(s) \approx \theta^T \phi(s)$, where $\phi(s)$ is the features of s .

* Using a model or simulator (assume to be discretized)

Informally, a simulator takes state s_t (continuous-valued) and action a_t as input, and outputs next-state s_{t+1} sampled according to $P_{s|at}$. One way to derive such a model is using a physics simulation. For example the inverted pendulum problem setting obeys the general laws of physics, thus by using physics knowledge one can derive the position and orientation the cart/pole will be in at time $t+1$ given state and action at time t .

An alternative approach would be to derive model from data. Suppose we execute n trials with each trial consisting T timesteps, we would then observe n state sequences:

$$S_0^{(1)}, a_0^{(1)} \rightarrow S_1^{(1)}, a_1^{(1)} \rightarrow S_2^{(1)}, a_2^{(1)} \rightarrow S_3^{(1)}, \dots \xrightarrow{a_T^{(1)}} S_T^{(1)}$$
$$S_0^{(n)}, a_0^{(n)} \rightarrow S_1^{(n)}, a_1^{(n)} \rightarrow \dots \rightarrow S_{T-1}^{(n)}, a_{T-1}^{(n)} \rightarrow S_T^{(n)}$$

We can then apply supervised learning algorithm to predict s_{t+1} as a function of s_t and a_t . For example in our helicopter problem setting, when helicopter is generally flying low-speed and experience no sharp turns, a linear model would work well, say we choose to fit

$$S_{t+1} = A s_t + B a_t, \text{ where } A \in \mathbb{R}^{d \times d}, B \in \mathbb{R}^{d \times d \text{ action}}$$

and further performing maximum likelihood estimate leads us to pick A, B by

$$\arg \min_{A, B} \sum_{t=0}^{T-1} \sum_{i=1}^n \|S_{t+1}^{(i)} - (A s_t^{(i)} + B a_t^{(i)})\|_2^2$$

Depending on our choice of features, modeling $S_{t+1} = A \phi(s_t) + B \phi'(a_t)$ can work better than our original linear model.

Having learned A, B , one option is to build a deterministic model, in which given input s_t and a_t , output S_{t+1} is exactly determined. Deterministic model works well for simulators, but once taken to the real-world, would usually work poorly on physical (robots, helicopters, etc). Alternatively we may also build a stochastic model, in which $S_{t+1} = As_t + Bs_t + \epsilon_t$, and noise term ϵ_t is usually modeled as $\epsilon_t \sim \mathcal{N}(0, \Sigma)$, where covariance matrix Σ can be estimated from data. **The exact choice of noise distribution matters less than the fact that we remember to add some noise.**

* Fitted Value Iteration

Recall our goal here is to approximate $V(s)$ as a linear or non-linear function of s , here we will use a supervised learning algorithm (linear regression in this case) and fit $V(s) = \theta^T \phi(s)$. Since value function on LHS suggest how well our model would be doing if start at state s (expected total payoff), it would be reasonable to choose features ($\phi(s)$) that help convey/evaluate the performance of model. For example in inverted pendulum problem, if the pole is leaning towards right, then we would want the car below to have larger velocity towards right, which suggest our feature function to contain \dot{x} as a entry, thus a reasonable feature function in this problem setting would be $\phi(s) = \begin{bmatrix} x^2 \\ \dot{x} \\ x \cdot \dot{x} \\ 0 \cdot \dot{x} \end{bmatrix}$. One potentially nice thing about model based RL is we can generate tons of data by using our model, thus we can afford to choose a lot of features (even irrelevant ones) without worrying about overfitting.

We will assume the problem has a continuous state space $S \in \mathbb{R}^d$, but that the action space is small and discrete, which is the case for most MDPs. Recall that in value iteration:

$$V(s) := R(s) + \gamma \max_a \int_{s'} P_{sa}(s') V(s') ds' = \max_a \mathbb{E}_{s' \sim P_{sa}} [R(s) + \gamma V(s')]$$

* Algorithm: Fitted Value Iteration

1. Sample $\{s^{(1)}, s^{(2)}, \dots, s^{(n)}\} \subseteq S$ randomly

2. Initialize $\theta := 0$. // Hence by $V(s) = \theta^T \phi(s)$, V would start with all entries being 0.

3. Repeat {

start for i
For $i = 1, 2, \dots, n$ { // having generated training sample $s^{(i)}$, what to estimate $V(s^{(i)})$

start for action $a \in A$ {

this corresponds to using a stochastic model, if deterministic
Sample $s'_1, s'_2, \dots, s'_K \sim P_{s^{(i)}a}$ (using model of MDP) set $K := 1$.

Set $g(a) := \mathbb{E}_{s' \sim P_{s^{(i)}a}} [R(s') + \gamma V(s')] = \frac{1}{K} \sum_{j=1}^K [R(s'^{(j)}) + \gamma V(s'_j)]$

end for actions.

Set $y^{(i)} = \max_a g(a)$. // Hence $y^{(i)}$ is an estimate of $R(s^{(i)}) + \gamma \max_{s' \sim P_{sa}} [V(s')]$.

end forⁱ } // Now $y^{(i)}$ is an estimate of $V(s^{(i)})$ which we want to model as $\theta^T \phi(s^{(i)})$

Set $\theta := \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (\theta^T \phi(s^{(i)}) - y^{(i)})^2$. // minimizing loss function of linear regression

} end

can plug in other supervised learning algos as well.

(or even neural network)

Unlike value iteration over discrete states, fitted value iteration cannot be proved to always converge

However in practice it often does converge (or approximately converge).

Finally, fitted value iteration outputs V which is an approximation of V^* , this implicitly defines our policy. When in state s , we would like to act according to $\pi^*(s) = \arg \max_a \mathbb{E}_{s' \sim P_{sa}} [V(s')]$, which in our approximation is $\arg \max_a \mathbb{E}_{s' \sim P_{sa}} [V(s')]$. Recall this is the inner loop of fitted value iteration.

but for systems that need frequent update, the reiterate sampling process can be computationally expensive. Say our model is $S_{t+1} = A s_t + B a_t + \epsilon_t$, during training time, adding noise is vital for robustness of policy, but when deploying in a physical simulator, one reasonable way is to set

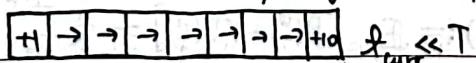
$\epsilon_t = 0$ and correspondingly $k_1 = 1$, which lead to optimal policy being $\arg \max_a V(f(s, a))$ [$S_{t+1} = f(s_t, a_t) + \epsilon_t$]

Conclusively, use a stochastic model during training phase and deterministic model during physical deployment.

• Finite Horizon MDP

We previously discussed the case where $\pi^*(s)$ is a stationary policy, that is, no matter the time stamp, our policy remains the same. In this section we put forward the case where $\pi^*(s)$ is a non-stationary policy, which, denotes as $\pi_t^*(s)$, varies over time.

We now assign a time duration for our MDP, and MDP becomes $(S, A, \{P_{sa}\}, T, R)$ with T being horizon time. The discount factor is not a must in this setting since our expected total payoff will always be a finite sum: $\mathbb{E}[R(s_0, a_0) + R(s_1, a_1) + \dots + R(s_T, a_T)]$. We now explicitly write out the action term in our reward function, this state-action reward can be used when we want to assign different reward to different actions. (e.g. For robot in a maze, staying still is preferred over wandering around, so we can charge this action; for helicopter setting, moving control sticks too aggressively is an action we want to penalize). One interesting property of finite horizon MDP is optimal policy depend on the time left, one dynamic illustration would be the policy derived for the following examples:



Let's focus on the finite horizon property and forget about state transition probability varying over time, and our inductive equation becomes: $V_t^*(s) = \max_a [R(s,a) + \sum_{s'} P_{sa}(s') V_{t+1}^*(s')]$ with $V_T^*(s)$ being $V_T^*(s) = \mathbb{E}[R(s_t, a_t) + \dots + R(s_T, a_T)]$ and correspondingly $\pi_t^*(s) = \arg \max_a [R(s,a) + \sum_{s'} P_{sa}(s') V_{t+1}^*(s')]$. One important thing is define $V_T^*(s) = \max_a R(s,a)$, which serves as base case for the following dynamic programming algorithm:

1. Calculate $V_T^*(s) = \max_a R(s,a)$, for $\forall s \in S$.

2. Use inductive step $V_t^*(s) = \max_a [R(s,a) + \sum_{s'} P_{sa}(s') V_{t+1}^*(s')]$ to sequentially compute $V_{T-1}^*, V_{T-2}^*, \dots$

This algorithm is value iteration in finite horizon setting.

Linear Quadratic Regulation

(with relatively small size)

In small set of problems with finite horizon property, linear quadratic Regulation is an efficient algorithm that is capable for the derivation of great control policies. To specify problem setting, say our MDP consists of five tuples $(S, A, \{P_{sa}\}, T, R)$, where $S = \mathbb{R}^n$, $A = \mathbb{R}^d$, and states evolving according to $S_{t+1} = AS_t + B a_t + w_t$, with $w_t \sim \mathcal{N}(0, \Sigma_w)$. One final assumption for applying LQR is that our reward function takes the form $R(s,a) = -(s^\top U s + a^\top V a)$ where $U \in \mathbb{R}^{n \times n}$, $V \in \mathbb{R}^{d \times d}$ and U, V are both PSD. One example for this is suppose that we want our helicopter to stay in state $s \approx 0$, then we can choose $U = I$ and $V = I$ to obtain $R(s,a) = -(||s||^2 + ||a||^2)$ which intuitively, penalizes drastic changes in state and action.

Now let's focus on calculating A, B . Previously in our discussion of model-based RL, we flew our helicopter m times to obtain m trajectories each consisting T time points, then we performed linear regression leading to minimize $\sum_{m=1}^M \sum_{t=0}^{T-1} ||S_{t+1} - (AS_t + B a_t)||^2$. LQR derives from the idea that if $S_{t+1} = f(S_t, a_t)$ and f is determined (e.g. in inverted pendulum f can be derived from laws of physics), can we linearize f around certain point^{*} so that A, B can be derived from the linear representation directly? This idea leads us to $S_{t+1} \approx f(\bar{s}_t, \bar{a}_t) + \nabla_s f(\bar{s}_t, \bar{a}_t)^T (\bar{s}_t - \bar{s}_t) + \nabla_a f(\bar{s}_t, \bar{a}_t)^T (\bar{a}_t - \bar{a}_t)$, where \bar{s}_t, \bar{a}_t are "typical points". We choose \bar{s}_t and \bar{a}_t to be the state/action where we want our model to spend most of its time on, since in a small vicinity of \bar{s}_t and \bar{a}_t , the linearized function is a fair estimate of the true curve. Based on the linearization, we estimate $A = \nabla_s f(\bar{s}_t, \bar{a}_t)^T$, $B = \nabla_a f(\bar{s}_t, \bar{a}_t)^T$, notice there are extra terms in the linear part, thus we adjust A to be $\tilde{A} = [f(\bar{s}_t, \bar{a}_t) - \nabla_s f(\bar{s}_t, \bar{a}_t)^T | A]$

and corresponding add an intercept term to s_t to obtain $s_t = \begin{bmatrix} \frac{1}{2} \\ s_t \end{bmatrix}$.

With linear quadratic regulation, hopefully we can model a MDP as a linear quadratic system, either s_{t+1} is a linear function of s_t, a_t , or $s_{t+1} = f(s_t, a_t)$ with f being non-linear mapping, but can be modelled linearly around the point of (\bar{s}_t, \bar{a}_t) . It turns out that with s_{t+1} being linear function of s_t, a_t and cost function being quadratic, we can compute V^* , which would be a quadratic function, exactly.

We will be developing a dynamic programming algorithm to demonstrate this point. Our base case would be $V_T^*(s_T) = \max_{a_T} R(s_T, a_T) = \max_{a_T} -(s_T^T U s_T + a_T^T V a_T)$, since $U, V \geq 0$, the best we can do to maximize $V_T^*(s_T)$ is to choose a_T s.t. $a_T^T V a_T = 0$, which leads to $V_T^*(s_T) = -s_T^T U s_T$, and $\pi_T^*(s_T) = \vec{0}$. We will assume that $V_{t+1}^*(s_{t+1}) = s_{t+1}^T \vec{\Phi}_{t+1} s_{t+1} + \psi_{t+1}$ where $\vec{\Phi}_{t+1} \in \mathbb{R}^{m \times m}$, $\psi_{t+1} \in \mathbb{R}$ is a quadratic function, we will then show after one iteration of dynamic programming, $V_t^*(s_t)$ will also be a quadratic function, which inductively proves our hypothesis (claim).

$$\text{Recall that } V_t^*(s_t) = \max_{a_t} (R(s_t, a_t) + \mathbb{E}_{\substack{s_{t+1} \sim P \\ a_{t+1}}} [V_{t+1}^*(s_{t+1})])$$

$$\begin{aligned} &= \max_{a_t} (-s_t^T U s_t + a_t^T V a_t + \mathbb{E}_{\substack{s_{t+1} \sim P \\ a_{t+1}}} [s_{t+1}^T \vec{\Phi}_{t+1} s_{t+1} + \psi_{t+1}]) \quad [\text{reminder: inside expectation is a scalar}] \\ &= \max_{a_t} \left[-(s_t^T U s_t + a_t^T V a_t) + \text{fr}(\vec{\Phi}_{t+1} \Sigma_w) + \mu_t^T \vec{\Phi}_{t+1} \mu_t + \psi_{t+1} \right], \text{ where } \mu_t = A s_t + B a_t. \\ &= \max_{a_t} \left[-(s_t^T U s_t + a_t^T V a_t) + s_t^T A^T \vec{\Phi}_{t+1} A s_t + \psi_{t+1} \right] + \left[-a_t^T V a_t + s_t^T A^T \vec{\Phi}_{t+1} B a_t + a_t^T B^T \vec{\Phi}_{t+1} A s_t + a_t^T B^T \vec{\Phi}_{t+1} B a_t \right] \\ &\quad \downarrow \quad \downarrow \quad \downarrow \quad \text{constant w.r.t } a_t \quad \text{big quadratic function of } a_t. \end{aligned}$$

If we take derivative of the above function and set to zero to solve a_t , we would obtain:

$$a_t = -(B^T \vec{\Phi}_{t+1} B - V)^{-1} B^T \vec{\Phi}_{t+1} A s_t$$

we will denote the linear part $-(B^T \vec{\Phi}_{t+1} B - V)^{-1} B^T \vec{\Phi}_{t+1} A$ as L_t , thus $\pi_t^*(s_t) = L_t s_t$. To summarize, in linear dynamical system with quadratic cost function, optimal action is a linear function of the state s_t . By plugging a_t back in we are able to obtain $V_t^*(s_t) = s_t^T [A^T (\vec{\Phi}_{t+1} + \vec{\Phi}_{t+1} B (V - B^T \vec{\Phi}_{t+1} B)^{-1} B^T \vec{\Phi}_{t+1}) A - U] s_t + \text{fr}(\vec{\Phi}_{t+1} \Sigma_w) + \psi_{t+1}$, which directly shows that $V_t^*(s_t)$ is also a quadratic function of s_t which supports our claim.

One thing to notice is that in $a_t = L_t s_t$, L_t is only represented through $\vec{\Phi}_{t+1}$ but not ψ_{t+1} , which indicates, when doing dynamic programming and iteratively update $\vec{\Phi}_{t+1}, \psi_{t+1}$ using $\vec{\Phi}_{t+1}$ and ψ_{t+1} , we can get rid of the calculation of ψ_{t+1} cause it's unrelated to a_t . Moreover, since Σ_w only appears in the calculation of ψ_{t+1} , then the specific choice of Σ_w won't affect our calculation.

of the optimal policy. To summarize, when modeling MDP as a linear dynamical system, the specific choice for the covariance of noise does not matter since it doesn't affect the optimal policy. Remember this is a property that is very specific to LQR, so do not overgeneralize it to other reinforcement learning algorithms.