

September 12, 2018 11:00 AM

create a function

$$\hat{y}(x) = y^{*2} \quad \text{prediction}$$

Error function choices

- Least squares:

- In 2-D: were given data $\{ \underset{\substack{\uparrow \\ \text{input}}}{x_i}, \underset{\substack{\uparrow \\ \text{output}}}{y_i} \}_{i=1}^N$

$$\hat{E}(\omega) = \sum_{i=1}^n \left(\underbrace{y_i}_{\text{observation}} - \underbrace{x_i \cdot \omega}_{\text{prediction}} \right)^2$$

b is bias can be noted.

find difference between
Observed output and input.

$$= \|y - xw\|_2^2 \leftarrow 2^{nd} \text{ norm}$$

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Observed output and input.
find difference. Motivation
of cost function

In multiple dimension, we have

$$(\vec{x}_i, y_i), \quad x \in \mathbb{R}^d$$

$$E(\vec{w}) = \|\vec{y} - X\vec{w}\|_2^2 \quad \text{2nd norm} \quad , \quad X = \begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix} \quad , \quad \vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

expanding the dot product
of itself

$$= (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w})$$

$$= \underbrace{\vec{w}^T X^T X \vec{w} - 2\vec{y}^T X \vec{w} + \vec{y}^T \vec{y}}_{\text{quadratic}}$$

Optimize it: Goal to minimize it. Take gradient and
set it to zero, (0). This will guarantee minimization

Linear Review

$$1. (AB)^T = B^T A^T$$

$$2. (AB)^{-1} = B^{-1} A^{-1} \leftarrow \text{assuming inverse exist}$$

$$3. (A^{-1})^T = (A^T)^{-1}$$

$$4. |AB| = |A||B|$$

↑
determinant(AB)

$$5. |A^{-1}| = \frac{1}{|A|}$$

6. A square matrix is orthogonal if every column vector

is orthogonal (checked product = I) and normalized

dot product of
vector by itself
is 1 ($\vec{x}^T \vec{x} = 1$)

Assume matrix is symmetric

Understand what orthogonal, symmetric,
and invertible mean.

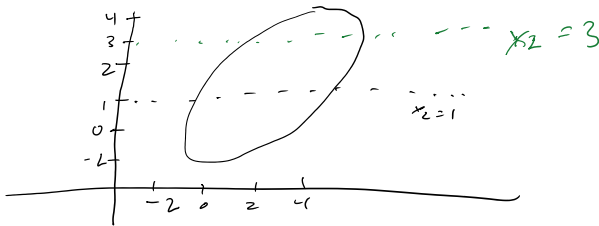
7. A square matrix "A" is not singular (invertible) if

Conditional (2D-gaussian)

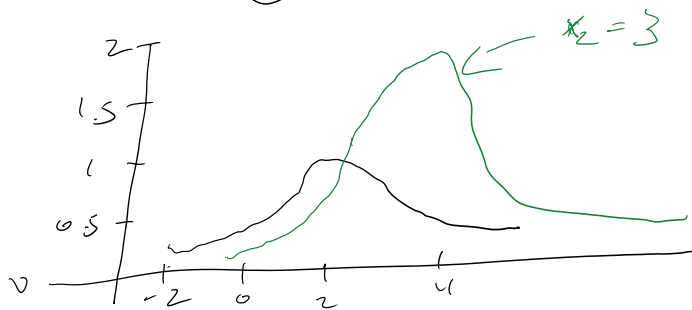
$$\text{Let } \vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, C = \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}, \Lambda = C^{-1} = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix}$$

\therefore the conditional distribution of x_1 given x_2 satisfies $x_1 | x_2 \sim \mathcal{N}(\mu_{x_1|x_2}, \Lambda_{11}^{-1})$, where where $\mu_{x_1|x_2} = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (x_2 - \mu_2)$, note that Λ_{11}^{-1} is not simply C_{11}

For ex



||



Diagonalization

Given a covariance matrix with eigen vectors \vec{u}_1, \vec{u}_2 and corresponding eigenvalues λ_1, λ_2

$$\text{let } U = \begin{bmatrix} | & | \\ \vec{u}_1 & \vec{u}_2 \\ | & | \end{bmatrix}, S = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix},$$

$$\text{then } CU = US, \quad C\vec{u}_1 = \lambda_1 S, \\ C\vec{u}_2 = \lambda_2 S,$$

since U is orthogonal, orthogonal if you

$$C = USU^{-1} = USU^T$$

$$\Rightarrow C^{-1} = (USU^T)^{-1} = U^{-1}S^{-1}U^{-T}$$

$$= US^{-1}U^T$$

Take the inverse
you get the transpose

$$\left(-\frac{1}{2}(\vec{x}-\vec{\mu})^T C^{-1}(\vec{x}-\vec{\mu})\right), \text{ sub } C^{-1} \text{ in,}$$

$$-\frac{1}{2}(\vec{x}-\vec{\mu})^T \underbrace{US^{-1}U^T}_{\text{change of variable}}(\vec{x}-\vec{\mu})$$

$$y = U^T(\vec{x}-\vec{\mu}) \Rightarrow \left(-\frac{1}{2}\vec{y}^T S^{-1}\vec{y}\right)$$

y is still a gaussian, with mean = 0, and covariance of S .

So S is diagonal matrix $\Rightarrow y_1, y_2$ are \perp

if covariance is diagonal then $y_1 \perp y_2$

\Rightarrow Eigen values of C are the variances along the principle directions given by the eigenvectors used for the diagonalization.

Positive Definite

A $d \times d$ matrix A is positive definite
iff $\forall \vec{z} \in \mathbb{R}^d - \{0\}, \vec{z}^T A \vec{z} > 0$

if $\vec{z}^T A \vec{z} \geq 0$, then positive semi-definite

why is $C > 0$?

1) Covariance is the second moment of a PDF (shifted by $\vec{\mu}$)

$$\Rightarrow C = E((\bar{x} - \bar{\mu})(\bar{x} - \bar{\mu})^T), \quad E(a^T)$$

If C is full rank (invertible) then

$$\bar{z}^T C \bar{z} = \bar{z}^T E((\bar{x} - \bar{\mu})(\bar{x} - \bar{\mu})^T) \bar{z}$$

$$= E(\bar{z}^T (\bar{x} - \bar{\mu})(\bar{x} - \bar{\mu})^T \bar{z})$$

$$= E(u^2), \quad u = (\bar{x} - \bar{\mu})^T \bar{z}$$

$$> 0$$

mean will always be positive because of the square.

$\Rightarrow C$ is positive definite.
if $C > 0$, eigen > 0

Why?

$$C \bar{z} = \lambda \bar{z} \quad \text{unit norm} \quad \bar{z} \text{ eigen vector}$$

$$\Leftrightarrow \bar{z}^T C \bar{z} = \lambda > 0 \quad \lambda \text{ eigen value}$$

$$\therefore \lambda > 0$$

Block Matrices (Reference: matrix identity)

A block matrix is a matrix that is interpreted as having been breaking into sections called blocks.

$$U = \begin{bmatrix} \boxed{u_1} & \boxed{u_2} \end{bmatrix}$$

$2 \times 1 \quad \quad 2 \times 1$

Useful in higher dimensions, use it to marginalize and condition

$$p(x_1, x_2)$$

$$p(x_1) = \int p(x_1, x_2) dx_2$$

$$p(x_2) = \int p(x_2, x_1) dx_1$$

For the case of Gaussians
Covariance

$$\Sigma = \begin{bmatrix} A & B \\ C & D \end{bmatrix}, \quad A \in \mathbb{R}^{m \times m} \\ D \in \mathbb{R}^{n \times n}$$

$$B = C^T$$

$$\text{if } \Sigma = \begin{bmatrix} A & 0 \\ 0 & D \end{bmatrix} \quad \text{then } \Sigma^{-1} = \begin{bmatrix} A^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix}$$

$$\det = \det(A) \det(D)$$

Regression - Quiz

helpful tips:

- # 1 generate noiseless data set
- # create own data set/training set
- # penalty for overfitting.

↳ avoid 1st second

Bayes' Rule

$$P(\bar{w} | D) = \underbrace{P(D | \bar{w})}_{\text{likelihood}} \times \underbrace{\frac{P(\bar{w})}{P(D)}}_{\text{evidence}} \quad \text{prior}$$

$$= \int P(D, w) dw = \int P(D | w) P(w) dw$$

\bar{w} parameter of given model

D = training data.

$$\rightarrow P(\bar{w} | D, M) = \frac{P(D | \bar{w}, M) \times P(\bar{w} | M)}{P(D | M)}$$

We care about different models.

In a, we are trial and erroring. But there is a more rigorous way using Bayes' Rule.

$$P(M | D) = \frac{P(D | M) P(M)}{P(D)} \leftarrow \text{constant}$$

→ Maximize posterior (give the "best" model)

Ex/ Estimating Gaussian Distribution

Suppose we are learning a gaussian distribution from n training data $\{\bar{x}_i\}_{i=1}^n$ and we want to know the best parameter $(\bar{\mu}, \Sigma)$ for this distribution.

$$\text{We already have the likelihood} = P(\bar{x}_{i=1:n} | \bar{\mu}, \Sigma)$$

$$= P(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n | \bar{\mu}, \Sigma)$$

$$= \prod_{i=1}^n P(\bar{x}_i | \bar{\mu}, \Sigma)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{(2\pi)^D}} \exp\left(-\frac{1}{2} (\bar{x}_i - \bar{\mu})^T \Sigma^{-1} (\bar{x}_i - \bar{\mu})\right) \quad \text{note } (2\pi)^D \leftarrow D \text{ is Dimension}$$

⇒ Because we restricted M = Gaussian, we

only need to maximize the likelihood (MLE).

⇒ Maximizing the above is complicated so instead, take the negative log-likelihood, minimizing it

$$\begin{aligned}\rightarrow L(\bar{\mu}, \bar{\Sigma}) &= -\ln(p(\bar{x}_{1:n} | \bar{\mu}, \bar{\Sigma})) \\ &= -\sum_{i=1}^N \ln p(\bar{x}_i | \bar{\mu}, \bar{\Sigma}) \\ &= \sum_{i=1}^N \left[\frac{(\bar{x}_i - \bar{\mu})^T \bar{\Sigma}^{-1} (\bar{x}_i - \bar{\mu})}{2} + \frac{N}{2} \ln |\bar{\Sigma}| + \frac{ND}{2} \ln(2\pi) \right]\end{aligned}$$

$$\bar{\mu}^* = \underset{\bar{\mu}}{\operatorname{argmax}} p(\bar{x}_{1:n} | \bar{\mu}, \bar{\Sigma}) = \underset{\bar{\mu}}{\operatorname{argmin}} L(\bar{\mu}, \bar{\Sigma})$$

$$\bar{\Sigma}^* = \underset{\bar{\Sigma}}{\operatorname{argmax}} p(\bar{x}_{1:n} | \bar{\mu}, \bar{\Sigma}) = \underset{\bar{\Sigma}}{\operatorname{argmin}} L(\bar{\mu}, \bar{\Sigma})$$

can solve them by

$$\frac{\partial L}{\partial \bar{\mu}} = 0 \quad \frac{\partial L}{\partial \bar{\Sigma}} = 0$$

$$\bar{\mu}^* = \frac{1}{N} \sum_{i=1}^N \bar{x}_i$$

$$\bar{\Sigma}^* = \frac{1}{N} (\bar{x}_i - \bar{\mu}^*) (\bar{x}_i - \bar{\mu}^*)^T$$

Entropy & Information Theory

- Entropy measures uncertainty of a distribution.
- Entropy is a key measure of uncertainty associated with a r.v.
- Use entropy for decision tree, k -L divergence, cross entropy
 - common measure of distance or MLE.
- Give a discrete r.v. y , which takes on k values, entropy is defined as $h = -\underbrace{E_{p(c)}[\log p(c)]}_{\text{expectation}} = -\sum_{c=1}^k p(c) \log(p(c))$

say we have $K=8$, $p(c=i) = \frac{1}{8} \quad \forall i=1, \dots, 8$

$$H = -8 \left(\frac{1}{8} \log \left(\frac{1}{8} \right) \right) = 3 \quad \uparrow \text{minimum bound}$$

This is used for data compression too.

1 hour, 4-5 questions. up to classification.

Short questions

- ↳ understand formulae least square
 - ↳ understand entropy
 - ↳ decision tree
 - ↳ gcc, linear regression
 - ↳ map and MLE
- } short answer

Understand Definition.

Optimization

We have $E(w)$, Find $\underset{w}{\operatorname{argmin}} E(w)$ ↙ weight that will minimize error

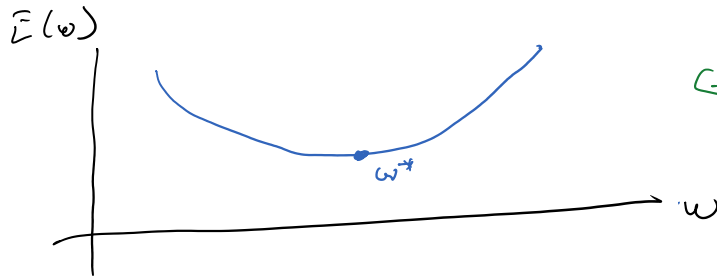
To find $\underset{w}{\operatorname{argmin}} E(w)$, take the gradient and isolate for w .

L.S regression - can get closed form solution

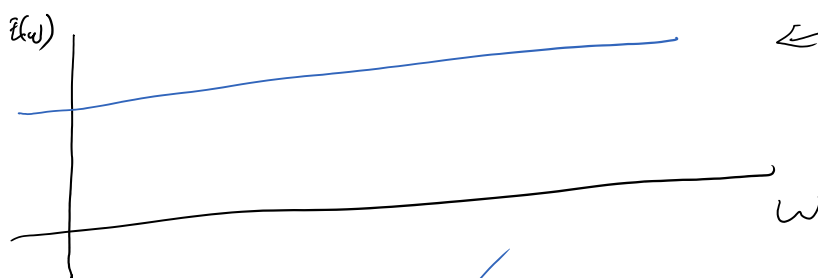
$$w^* = \underbrace{(X^T X)^{-1}}_{\substack{\text{computation} \\ \text{for inverse is} \\ \text{costly}}} X^T y$$

L.S regression cost function

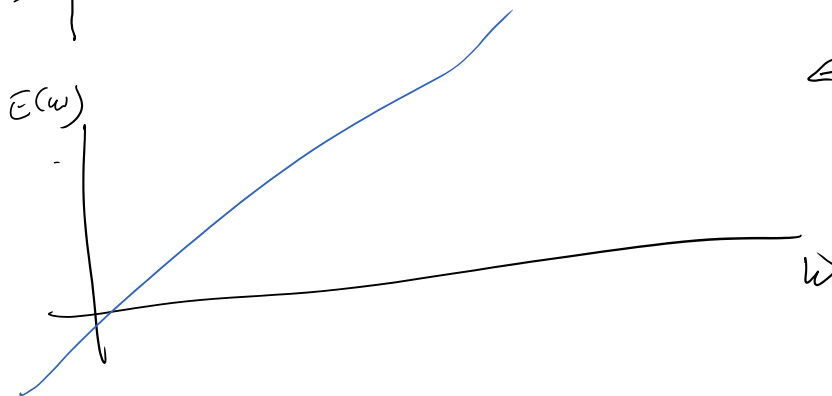
$$E(w) = w^T X^T X w - 2 y^T X w + y^T y$$



← draw in terms of parameter space



← no point, cause you won't be able to learn



← no minimum
can go negative infinity.

L2 loss is popular because you can derive it easily.

Gradient Descent (uses first order Taylor expansion)

- Can approximate error function

$$\approx E(\bar{w}_0) + (\bar{w} - \bar{w}_0) \frac{\partial E(\bar{w}_0)}{\partial \bar{w}}$$

- Goal: to minimize error

$$\Rightarrow \min E(\bar{w})$$

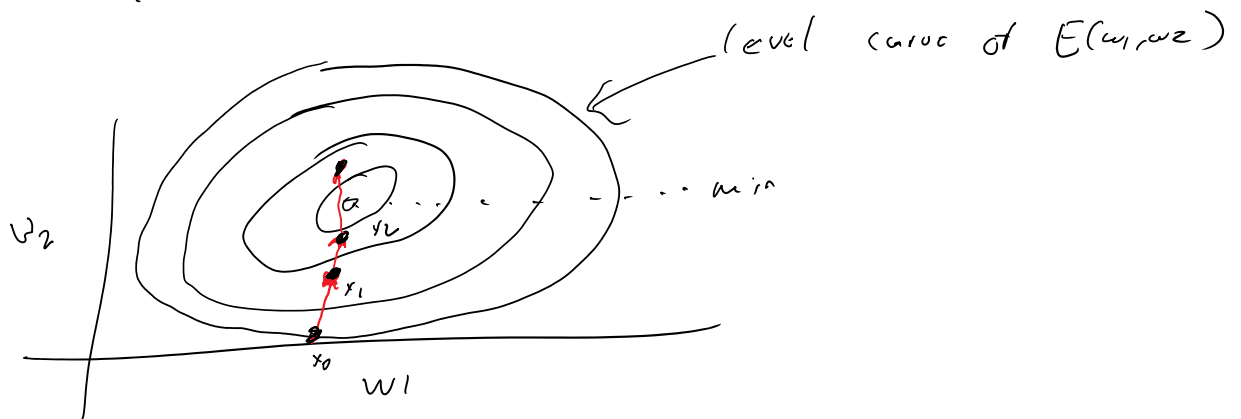
$$\Rightarrow \min_{\bar{h}} E(\bar{w}_0 + \bar{h}) \approx \min_{\bar{h}} E(\bar{w}_0) + \bar{h} \frac{\partial E(\bar{w}_0)}{\partial \bar{h}}$$

$$\Rightarrow \min_{\bar{h}} \frac{\partial E(\bar{w}_0)}{\partial \bar{h}} \leftarrow \text{take gradient}$$

$$\Rightarrow \bar{h}^* = -\nabla E(\bar{w}_0)$$

$$\Rightarrow w_{i+1} = w_i - \underset{\substack{\uparrow \\ \text{step size}}}{\alpha} \frac{\partial E(\bar{w}_i)}{\partial \bar{w}} \quad \left. \vphantom{\frac{\partial E(\bar{w}_i)}{\partial \bar{w}}} \right\} \text{optimization step}$$

$\bar{w}_0 = ? \rightarrow$ anything except 0



Gradient decent is used to get local minimum.

set step size as a hyper parameter