Administrivia

- Ed!!!
- Everything due Tuesdays, midnight
- Worksheets, HW0, HW1
 - We'll be lax about the deadlines for weeks 0,1
 - but do everything
- Surveys: DO THEM!!!!
- Office hours:wiki: "people/office hours"
- Friday review now in Wu and Chen

To study for 5200, Learn the concepts!

- Each day, note repeated terms
 - norm, entropy, log-likelihood, MLE, MAP ...
- What is the mathematical definition?
 - And the intuitive English description?
- Where is it used?
 - What is an alternative to it?
- What is its behavior in different limits?
 - N infinite, p goes to zero or 1, ...

Linear Regression

Lyle Ungar

Learning objectives Be able to derive MLE & MAP regression and the associated loss functions Recognize *scale invariance*

Slide 3

MLE estimates

A) $\operatorname{argmax}_{\theta} p(\theta|\mathbf{D})$ B) $\operatorname{argmax}_{\theta} p(\mathbf{D}|\theta)$ C) $\operatorname{argmax}_{\theta} p(\mathbf{D}|\theta)p(\theta)$ D) None of the above

A, B, C or D				
A				
В				
С				
D				

MAP estimates

A) $\operatorname{argmax}_{\theta} p(\theta|\mathbf{D})$ B) $\operatorname{argmax}_{\theta} p(\mathbf{D}|\theta)$ C) $\operatorname{argmax}_{\theta} p(\mathbf{D}|\theta)p(\theta)$ D) None of the above

A, B, C or D				
٨				
R				
C				
D				
U				

MLE vs MAP

We will use:
A) MLE
B) MAP



• Why?

Gaussian MAP

When does this become the MLE?

$$\hat{\mu}_{MAP} = rac{rac{\mu_0}{\sigma_0^2} + \sum_i rac{x_i}{\sigma^2}}{rac{1}{\sigma_0^2} + rac{n}{\sigma^2}}$$

Based on prior

$$p(\mu|\mu_0,\sigma_0) = rac{1}{\sigma_0\sqrt{2\pi}}e^{rac{-(\mu-\mu_0)^2}{2\sigma_0^2}}$$

Consistent estimator

A consistent estimator (or asymptotically consistent estimator) is an estimator — a rule for computing estimates of a parameter θ — having the property that as the number of data points used increases indefinitely, the resulting sequence of estimates converges in probability to the true parameter θ.

https://en.wikipedia.org/wiki/Consistent_estimator

Which is consistent for our coinflipping example?



Slide 9

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <u>http://www.cs.cmu.edu/~awm/tutorials</u>. Comments and corrections gratefully received.

An introduction to regression

mostly by Andrew W. Moore but with many modifications by Lyle Ungar

Two interpretations of regression

- Linear regression
 - $\hat{y} = \mathbf{W} \cdot \mathbf{x}$
- Probabilistic/Bayesian (MLE and MAP)
 - $y(\mathbf{x}) \sim N(\mathbf{w} \cdot \mathbf{x}, \sigma^2)$
 - MLE: argmax_w p(**D**|w)
 - MAP: argmax_w p(**D**|w)p(w)
- Error minimization
 - $\|\boldsymbol{y} \boldsymbol{X}\boldsymbol{w}\|_p^p + \lambda \|\boldsymbol{w}\|_q^q$

here: argmax_w p(y|w,X)

Single-Parameter Linear Regression

Linear Regression



Linear regression assumes that the expected value of the output given an input, E[y|x], is linear in x.

Simplest case: $\hat{y}(x) = wx$ for some unknown w. (x,w scalars)

Given the data, we can estimate w.

One parameter linear regression

Assume that the data is formed by

 $y_i = wx_i + noise_i$

where...

- *noise*_{*i*} is independent N(0, σ^2)
- variance σ^2 is unknown
- y(x) then has a normal distribution with
- mean wx
- variance σ^2

Copyright $\ensuremath{\textcircled{C}}$ 2001, 2003, Andrew W. Moore

Bayesian Linear Regression

p(y|w,x) is Normal(mean: *wx*, variance: σ^2) y ~ N(wx, σ^2)

We have a data $(x_1, y_1) (x_2, y_2) \dots (x_n, y_n)$

We want to estimate *w* from the data.

 $p(w|x_1, x_2, x_3, ..., x_n, y_1, y_2..., y_n) = P(w|D)$

•You can use BAYES rule to find a posterior distribution for *w* given the data.

•Or you could do Maximum Likelihood Estimation

Maximum likelihood estimation of w

MLE asks :

"For which value of w is this data most likely to have happened?"

<=>

For what w is

 $p(y_1, y_2...y_n | w, x_1, x_2, x_3,...x_n) \text{ maximized?}$ $< = > \qquad \qquad \prod_{i=1}^{n} p(y_i | w, x_i) \text{ maximized?}$ For what w is

i=1



Result: MLE = L₂ error

• MLE with Gaussian noise is the same as minimizing the L₂ error

$$\operatorname{argmin}_{i=1}^{n} (y_i - w x_i)^2$$

Linear Regression

The maximum likelihood *w* is the one that minimizes sum-of-squares of residuals



 $r_i = y_i - W x_i$

 $E = \sum_{i} (y_{i} - wx_{i})^{2}$ $= \sum_{i} y_{i}^{2} - (2\sum_{i} x_{i}y_{i})w + (\sum_{i} x_{i}^{2})w^{2}$

We want to minimize a quadratic function of *w*.

Linear Regression

p(w

The sum of squares is minimized when

$$w = \frac{\sum x_i y_i}{\sum x_i^2}$$

The maximum likelihood model is

 $\hat{y}(x) = wx$

We can use it for prediction

Copyright $\ensuremath{\mathbb{C}}$ 2001, 2003, Andrew W. Moore

Note: In Bayesian stats you'd have ended up with a prob. distribution of *w*

W

And predictions would have given a prob. distribution of expected output

Often useful to know your confidence. Max likelihood can give some kinds of confidence, too.

But what about MAP?

• MLE

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



But what about MAP?

• MAP

$$\operatorname{argmax} \prod_{i=1}^{n} p(y_i | w, x_i) p(w)$$

- We assumed
 - $y_i \sim N(w x_i, \sigma^2)$
- Now add a prior assumption that
 - w ~ N(0, γ^2)

For what *w* is

$$\prod_{i=1}^{n} p(y_i | w, x_i) p(w) \text{ maximized}?$$
For what *w* is

$$\prod_{i=1}^{n} \exp(-\frac{1}{2}(\frac{y_i - wx_i}{\sigma})^2) \exp(-\frac{1}{2}(\frac{w}{\gamma})^2) \text{maximized}?$$
For what *w* is

For what *W* is

$$\sum_{i=1}^{n} -\frac{1}{2} \left(\frac{y_i - wx_i}{\sigma} \right)^2 - \frac{1}{2} \left(\frac{w}{\gamma} \right)^2 \text{ maximized}?$$

For what *w* is

$$\sum_{i=1}^{n} (y_i - wx_i)^2 + (\frac{\sigma w}{\gamma})^2 \text{ minimized}?$$

Ridge Regression is MAP

• MAP with a Gaussian prior on *w* is the same as minimizing the L₂ error plus an L₂ penalty on w

$$\operatorname{argmin} \sum_{i=1}^{n} (y_i - w x_i)^2 + \lambda w^2 \qquad \lambda = \sigma^2 / \gamma^2$$

- This is called
 - Ridge regression
 - Shrinkage
 - Tikhonov Regularization

Copyright C 2001, 2003, Andrew W. Moore

• $w = x^T y/(x^T x + \lambda)$ = $(x^T x + \lambda)^{-1} x^T y$

Multivariate Linear Regression

What if the inputs are vectors?



Write matrix X and Y thus:

$$\mathbf{X} = \begin{bmatrix} \dots \mathbf{X}_{1} \dots \mathbf{X}_{1} \\ \dots \mathbf{X}_{2} \dots \\ \vdots \\ \dots \mathbf{X}_{n} \dots \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_{1} \\ y_{2} \\ \vdots \\ y_{n} \end{bmatrix}$$

(*n* data points; Each input has *p* features) The linear regression model assumes $\hat{y}(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} = w_1 x_1 + w_2 x_2 + \dots w_p x_p$

The maximum likelihood estimate (MLE) is $\mathbf{w} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}(\mathbf{X}^{\mathsf{T}}\mathbf{y})$

X^T**X** is p x p **X**^T**y** is p x 1

The MAP estimate is $w = (X^TX + \lambda I)^{-1}(X^Ty)$

 $X^T X$ is $p \ge p$ I is the $p \ge p$ identity matrix $X^T y$ is $p \ge 1$

What about a constant term?

Linear data usually does not go through the origin.

Statisticians and Neural Net folks all agree on a simple obvious hack. Can you guess??



The constant term

• The trick: create a fake input " x_0 " that is always 1



Before: $Y=w_1X_1 + w_2X_2$...has to be a poor model X_0 Y X_1 X_2 1 2 16 4 1 3 17 4 1 5 20 5

After: $Y = w_0 X_0 + w_1 X_1 + w_2 X_2$ $= w_0 + w_1 X_1 + w_2 X_2$...has a fine constant term

L₁ regression

OLS = L₂ regression minimizes $p(y|w,x) \sim \exp(-||y-x \cdot w||_2^2/2\sigma^2) \longrightarrow \operatorname{argmin}_w ||y-Xw||_2^2$ L₁ regression: $p(y|w,x) \sim \exp(-||y-x \cdot w||_1/2\sigma^2) \longrightarrow \operatorname{argmin}_w ||y-Xw||_1$

Scale Invariance

- Rescaling does not affect decision trees or OLS
 - They are scale invariant
- Rescaling does affect Ridge regression
 - Because it preferentially shrinks 'large' coefficients

Scale invariance

 A method is scale invariant if multiplying any feature (any column of X) by a nonzero constant does not change the predictions made (after retraining)

y	X	У	X
1	0.2 3 5	1	0.2 3 10
0	0.3 4 5	0	0.3 4 10
0	0.1 4 4	0	0.1 4 8
1	0.2 2 2	1	0.2 2 4

K-NN is scale invariant?

Linear regression is scale invariant? Ridge regression is scale invariant?



What we have seen

- MLE with Gaussian noise minimizes the L₂ error
 - Other noise models will give other loss functions
- MAP with a Gaussian prior gives Ridge regression
 - Other priors will give different penalties
- Both are consistent
- Linear Regression is scale invariant; Ridge is not
- One can
 - do nonlinear regression
 - make nonlinear relations linear by transforming the features

 $//\boldsymbol{y} - \boldsymbol{X} \boldsymbol{w} / /_{p}^{p} + \lambda / / \boldsymbol{w} / /_{q}^{q}$



Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

Silae 37

What questions do you have on today's class?

Тор

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app



Heteroscedasticity…

Linear Regression with varying noise

Regression with varying noise

 Suppose you know the variance of the noise that was added to each datapoint.



MLE estimation with varying noise

argmax log $p(y_1, y_2, ..., y_R | x_1, x_2, ..., x_R, \sigma_1^2, \sigma_2^2, ..., \sigma_R^2, w) =$ Assuming independence \mathcal{W} $\underset{W}{\operatorname{argmin}} \sum_{i=1}^{R} \frac{(y_i - wx_i)^2}{\sigma_i^2} = \operatorname{among noise and then}_{\begin{array}{c} \text{plugging in equation for} \\ \begin{array}{c} \text{Gaussian and simplifying.} \end{array}}$ $\left(w \text{ such that } \sum_{i=1}^{R} \frac{x_i(y_i - wx_i)}{\sigma_i^2} = 0 \right) =$ Setting dLL/dw equal to zero Trivial algebra $\frac{\left(\sum_{i=1}^{R} \frac{x_i y_i}{\sigma_i^2}\right)}{\left(\sum_{i=1}^{R} \frac{x_i^2}{\sigma_i^2}\right)}$ Note that "R" here is what we call "n"

This is Weighted Regression

• We are minimizing the *weighted* sum of squares



where the weight for i'th datapoint is $\frac{1}{\sigma_{i}^{2}}$

Nonlinear Regression

Nonlinear Regression

• Suppose you know that y is related to a function of x in such a way that the predicted values have a non-linear dependence on w, e.g.:



Nonlinear MLE estimation

 $\operatorname{argmax} \log p(y_1, y_2, ..., y_R | x_1, x_2, ..., x_R, \sigma, w) =$ $\operatorname{argmin}_{W} \sum_{i=1}^{R} \left(y_i - \sqrt{w + x_i} \right)^2 =$ $\operatorname{Assuming i.i.d. and}_{\text{then plugging in}}_{\text{equation for Gaussian}}_{\text{and simplifying.}}$ $\left(w \operatorname{such that} \sum_{i=1}^{R} \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0 \right) =$ $\operatorname{Setting dLL/dw}_{\text{equal to zero}}$

Nonlinear MLE estimation

argmax log $p(y_1, y_2, ..., y_R | x_1, x_2, ..., x_R, \sigma, w) =$ Assuming i.i.d. and ${\mathcal W}$ then plugging in argmin $\sum_{i=1}^{R} (y_i - \sqrt{w + x_i})^2 =$ equation for Gaussian and simplifying. \mathcal{W} Setting dLL/dw $\left(w \text{ such that } \sum_{i=1}^{R} \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0\right) = 1$ equal to zero We' re down the algebraic toilet So guess what we do?

Nonlinear MLE estimation

