

Lecture 3:

Maximum Likelihood Estimation (MLE), Linear Regression, Linear Models



Markus Hohle

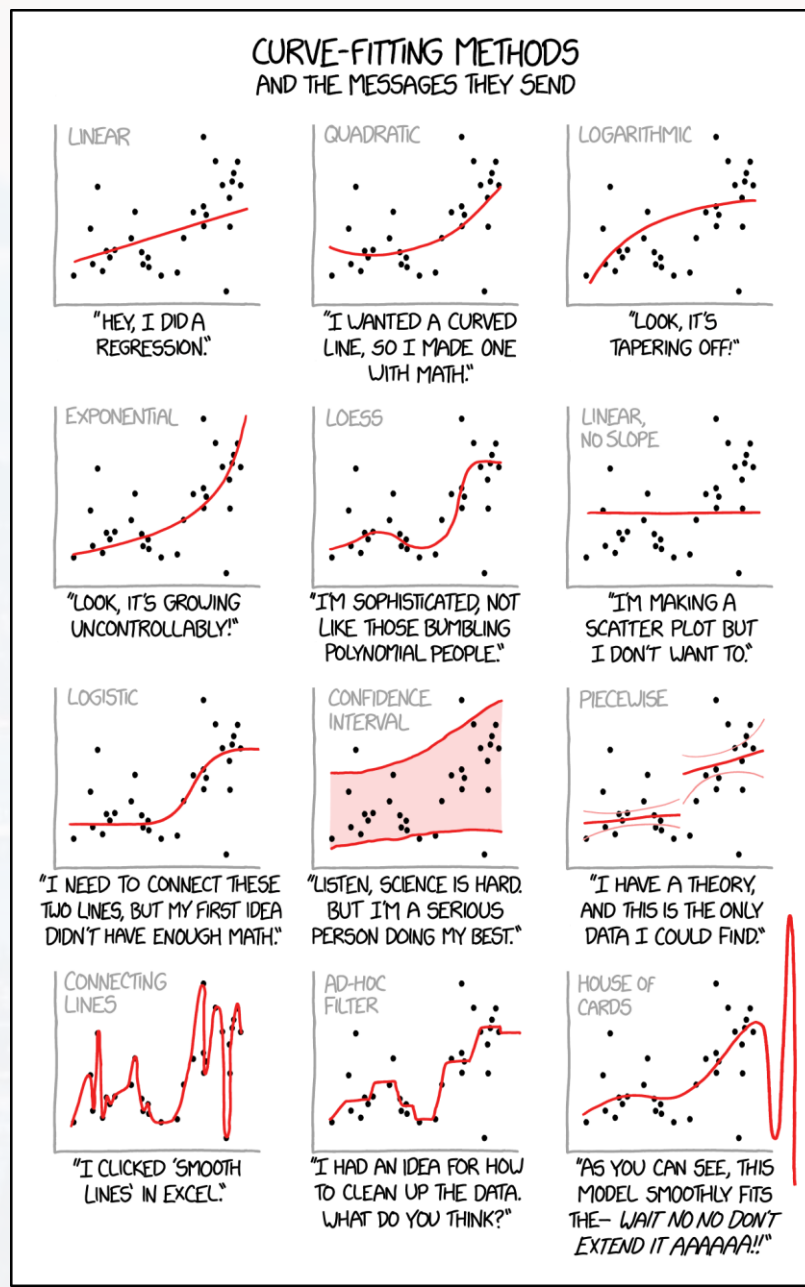
University California, Berkeley

Bayesian Data Analysis and
Machine Learning for Physical
Sciences



Course Map

Module 1	Maximum Entropy and Information, Bayes Theorem
Module 2	Naive Bayes, Bayesian Parameter Estimation, MAP
Module 3	MLE, Lin Regression, Model selection: Comparing Distributions
Module 4	Model Selection: Bayesian Signal Detection
Module 5	Variational Bayes, Expectation Maximization
Module 6	Stochastic Processes
Module 7	Monte Carlo Methods
Module 8	Markov Models, Graphs
Module 9	Machine Learning Overview, Supervised Methods
Module 10	Unsupervised Methods
Module 11	ANN: Perceptron, Backpropagation
Module 12	ANN: Basic Architecture, Regression vs Classification, Backpropagation again
Module 13	Convolution and Image Classification and Segmentation
Module 14	TBD (GNNs)
Module 15	TBD (RNNs and LSTMs)
Module 16	TBD (Transformer and LLMs)



Outline

Maximum Likelihood Estimation (MLE)

- idea
- examples

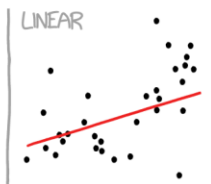
Linear Models

- classical model
- regularization
- extensions

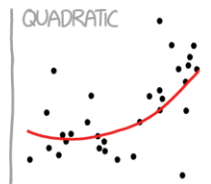
Regression



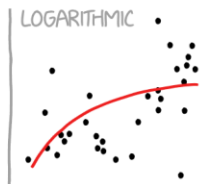
CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



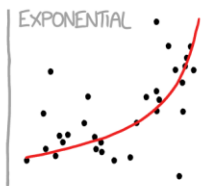
"HEY, I DID A
REGRESSION."



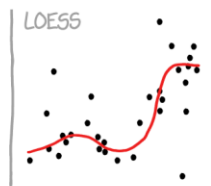
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



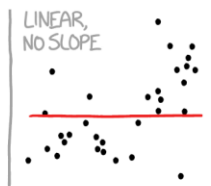
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



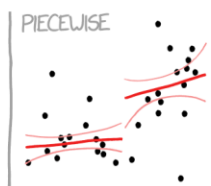
"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



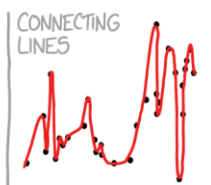
"I NEED TO CONNECT THESE
TWO LINES, BUT MY FIRST IDEA
DIDN'T HAVE ENOUGH MATH."



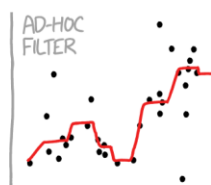
"LISTEN, SCIENCE IS HARD.
BUT I'M A SERIOUS
PERSON DOING MY BEST."



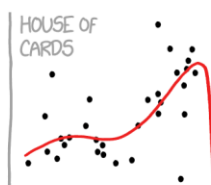
"I HAVE A THEORY,
AND THIS IS THE ONLY
DATA I COULD FIND."



"I CLICKED 'SMOOTH
LINES' IN EXCEL."



"I HAD AN IDEA FOR HOW
TO CLEAN UP THE DATA.
WHAT DO YOU THINK?"



"AS YOU CAN SEE, THIS
MODEL SMOOTHLY FITS THE-
WAIT NO NO DON'T
EXTEND IT AAAAAA!!"

Outline

Maximum Likelihood Estimation (MLE)

- idea
- examples

Linear Models

- classical model
- regularization
- extensions

Regression



Bayesian:

idea

examples

$$P(\theta|D) = \frac{\overset{\text{likelihood function}}{P(D|\theta)} \overset{\text{prior}}{P(\theta)}}{\underset{\text{evidence (const wrt } \theta \text{)}}{P(D)}} \quad \int P(\theta|D) d\theta = 1 \quad \text{parameter } \theta$$

“Classical” approach:

Maximum Likelihood Estimation

$x = \{x_1, x_2, \dots, x_i, \dots, x_N\}$ data set of N independent observations

$\theta = \{\theta_1, \theta_2, \dots, \theta_i, \dots, \theta_M\}$ set of parameters θ

$\theta = q$ (binomial)

$\theta = \{\mu, \sigma^2\}$ (normal dist)

$f_i(\theta, x_i)$

distribution f_i from which x_i has been drawn



“Classical” approach:

Maximum Likelihood Estimation

idea

examples

$x = \{x_1, x_2, \dots, x_i, \dots, x_N\}$ data set of N independent observations

$\theta = \{\theta_1, \theta_2, \dots, \theta_i, \dots, \theta_M\}$ set of parameters θ

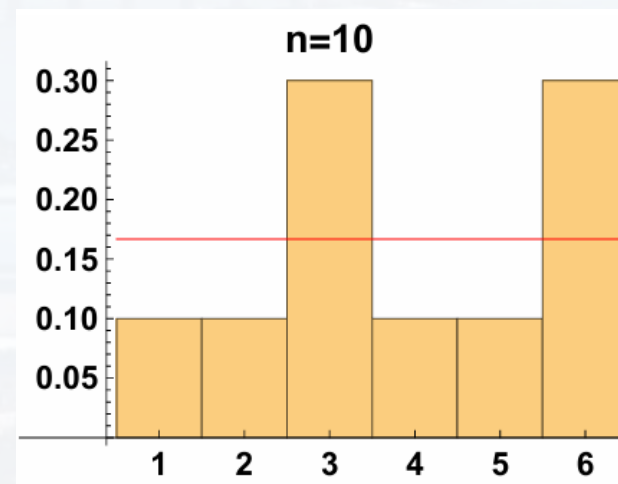
$f_i(\theta, x_i)$ distribution f_i from which x_i has been drawn

$$L(\theta, x) = \prod_{i=1}^N f_i(\theta, x_i) \quad \text{joint density}$$

goal: finding those $\hat{\theta}$ that maximize $L(\theta, x)$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \{L(\theta, x)\}$$

We aim to find the most likely one!



$\sigma = 2, \mu = 3.5$

$\sigma = 2, \mu = 5.0$

$\sigma = 1.5, \mu = 3.5$

$\sigma = 7.0, \mu = 1.0$

....and so on



“Classical” approach:

Maximum Likelihood Estimation

idea

examples

$$L(\theta, x) = \prod_{i=1}^N f_i(\theta, x_i) \quad \text{joint density}$$

goal: finding those $\hat{\theta}$ that maximize $L(\theta, x)$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \{L(\theta, x)\}$$

x :	data set
θ :	parameter
$f_i(\theta, x_i)$:	density function

note:

- often we use $l(\theta, x) = \ln[L(\theta, x)]$ for convenience

- we find $\hat{\theta}$ via $\frac{\partial L(\theta, x)}{\partial \theta_i} = 0$ for all θ_i

- equivalent to **MAP** assuming **an uniform prior**

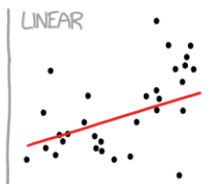
- evaluate, if extreme is indeed a maximum:

Hessian matrix $H = \left\{ \frac{\partial^2 l(\theta, x)}{\partial \theta_i \partial \theta_j} \right\}$ has to be **locally** (around $\hat{\theta}$) **concave**

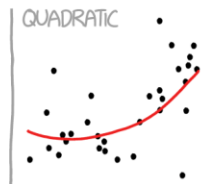
- we can add constraints using **Lagrangian Multipliers**



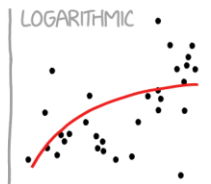
CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



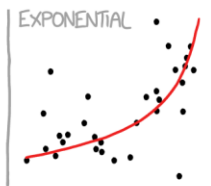
"HEY, I DID A
REGRESSION."



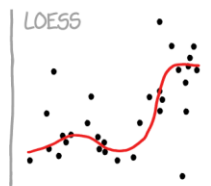
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



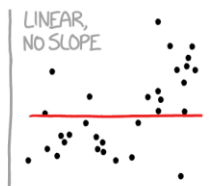
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



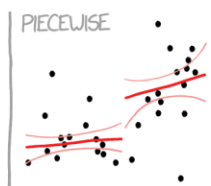
"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



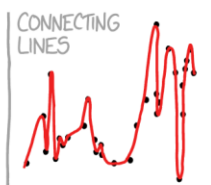
"I NEED TO CONNECT THESE
TWO LINES, BUT MY FIRST IDEA
DIDN'T HAVE ENOUGH MATH."



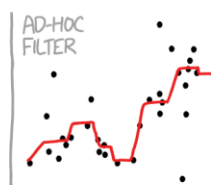
"LISTEN, SCIENCE IS HARD.
BUT I'M A SERIOUS
PERSON DOING MY BEST."



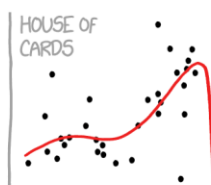
"I HAVE A THEORY,
AND THIS IS THE ONLY
DATA I COULD FIND."



"I CLICKED 'SMOOTH
LINES' IN EXCEL."



"I HAD AN IDEA FOR HOW
TO CLEAN UP THE DATA.
WHAT DO YOU THINK?"



"AS YOU CAN SEE, THIS
MODEL SMOOTHLY FITS THE-
WAIT NO NO DON'T
EXTEND IT AAAAAA!!"

Outline

Maximum Likelihood Estimation (MLE)

- idea
- examples

Linear Models

- classical model
- regularization
- extensions

Regression



idea
examples

example I: binomial distribution

$$x = \{HHTHTTT \dots\} \quad \theta = q$$

$$f(q|x) = \binom{n}{k} q^k (1-q)^{n-k} \quad \text{equivalent to MAP assuming an uniform prior} \quad \frac{P(D|\theta)P(\theta)}{P(D)}$$

$$\frac{df}{dq} = \binom{n}{k} [kq^{k-1}(1-q)^{n-k} - q^k(n-k)(1-q)^{n-k-1}] = 0$$

$$\text{ignoring } q = 1 \text{ and } q = 0: \quad k(1-q) = q(n-k) \quad q_{best} = \frac{k}{n}$$

$$q_0 = \frac{k}{n}$$

$$\sigma = \sqrt{\frac{q_0(1-q_0)}{n}}$$

$$q = \frac{k}{n} \pm \sqrt{\frac{q_0(1-q_0)}{n}}$$

last time:
Lagrange approximation from BPE



example II: normal distribution

idea

examples

$$x = \{x_1, x_2, \dots, x_i, \dots, x_N\} \quad \theta = \{\mu, \sigma^2\}$$

$$f_i(\theta, x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}}$$

$$L(\theta, x) = \prod_{i=1}^N f_i(\theta, x_i) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2}$$

$$l(\theta, x) = \ln[L(\theta, x)] = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

$$\frac{\partial l(\theta, x)}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^N (x_i - \mu)^2 = 0$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$



example II: normal distribution

idea

examples

$$x = \{x_1, x_2, \dots, x_i, \dots, x_N\} \quad \theta = \{\mu, \sigma^2\}$$

$$f_i(\theta, x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}}$$

$$L(\theta, x) = \prod_{i=1}^N f_i(\theta, x_i) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2}$$

$$l(\theta, x) = \ln[L(\theta, x)] = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

$$\frac{\partial l(\theta, x)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0$$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$



example II: normal distribution

idea

examples

$$x = \{x_1, x_2, \dots, x_i, \dots, x_N\} \quad \theta = \{\mu, \sigma^2\}$$

$$f_i(\theta, x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}}$$

$$L(\theta, x) = \prod_{i=1}^N f_i(\theta, x_i) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2}$$

$$l(\theta, x) = \ln[L(\theta, x)] = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

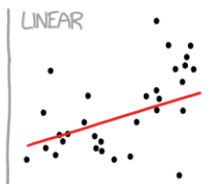
$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

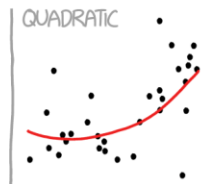
$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2 = \frac{1}{N} \sum_{i=1}^N \left(x_i - \frac{1}{N} \sum_{i=1}^N x_i \right)^2 = \frac{1}{N} \sum_{i=1}^N x_i^2 - \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N x_i x_j$$



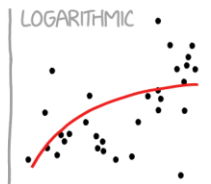
CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



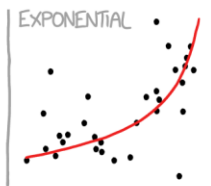
"HEY, I DID A
REGRESSION."



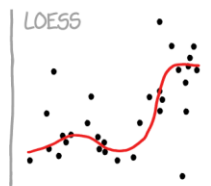
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



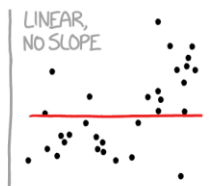
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



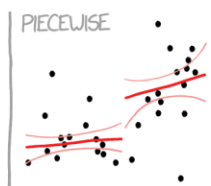
"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



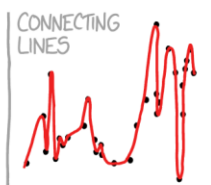
"I NEED TO CONNECT THESE
TWO LINES, BUT MY FIRST IDEA
DIDN'T HAVE ENOUGH MATH."



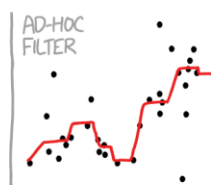
"LISTEN, SCIENCE IS HARD.
BUT I'M A SERIOUS
PERSON DOING MY BEST."



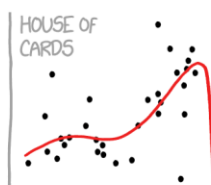
"I HAVE A THEORY,
AND THIS IS THE ONLY
DATA I COULD FIND."



"I CLICKED 'SMOOTH
LINES' IN EXCEL."



"I HAD AN IDEA FOR HOW
TO CLEAN UP THE DATA.
WHAT DO YOU THINK?"



"AS YOU CAN SEE, THIS
MODEL SMOOTHLY FITS THE-
WAIT NO NO DON'T
EXTEND IT AAAAAA!!"

Outline

Maximum Likelihood Estimation (MLE)

- idea
- examples

Linear Models

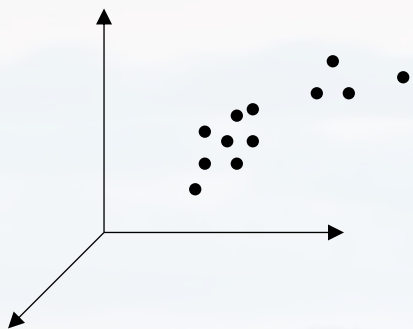
- classical model
- regularization
- extensions

Regression



idea: data point y_k in N dimensional feature space

classical model
regularizations
extensions



$$y_k = f(x_1, \dots, x_n, \dots, x_N) + \epsilon_k \text{ for each data point } k$$

$$y_k = \beta_0 + \sum_{n=1}^N \beta_n x_n + \epsilon_k$$

y:	response
x:	regressors (assumed to be independent)
β:	factors (how a regressor contributes to the response)
β_0:	intercept
ϵ:	error (stochasticity of the data, assumed to be normally dist.)

$$\underbrace{\begin{pmatrix} y_1 \\ \dots \\ y_k \\ \dots \\ y_K \end{pmatrix}}_Y = \underbrace{\begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} & \dots & x_{1N} \\ \dots & \dots & \dots & & \dots & & \dots \\ 1 & x_{k1} & & & x_{kn} & & \\ \dots & \dots & & & \dots & & \dots \\ 1 & \dots & & & \dots & & \dots \\ 1 & x_{K1} & x_{K2} & \dots & x_{Kn} & \dots & x_{KN} \end{pmatrix}}_X \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \dots \\ \beta_n \\ \dots \\ \beta_N \end{pmatrix}}_{\beta} + \underbrace{\begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \dots \\ \epsilon_k \\ \dots \\ \epsilon_K \end{pmatrix}}_{\epsilon}$$

$$Y = X\beta + \epsilon$$



$$Y = X\beta + \varepsilon$$

y :	response
x :	regressors (assumed to be independent)
β :	factors (how a regressor contributes to the response)
β_0 :	intercept
ε :	error (stochasticity of the data, assumed to be normally dist.)

classical model
regularizations
extensions

fitting: finding the best β in terms of minimizing the errors

$$(Y - X\beta)^T(Y - X\beta) = \sum_k \varepsilon_k^2$$

$$\frac{\partial}{\partial \beta} \sum_k \varepsilon_k^2 = 0$$

$$\beta_{best} = \hat{\beta} = (X^T X)^{-1} X^T Y$$

the model

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

X and Y are all observables

$$\text{hat matrix } H := X(X^T X)^{-1} X^T$$

$$H = H^T \quad (\text{symmetry})$$
$$HH = H \rightarrow H^n = H \quad (\text{idempotency})$$



$$Y = X\beta + \varepsilon$$

y : response
 x : regressors (assumed to be independent)
 β : factors (how a regressor contributes to the response)
 β_0 : intercept
 ε : error (stochasticity of the data, assumed to be normally dist.)

classical model
regularizations
extensions

evaluating the result:

$$\hat{\varepsilon} = Y - X\hat{\beta} = Y - \hat{Y} = (I - H)Y$$

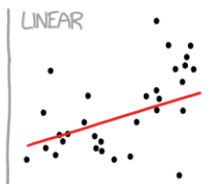
hat matrix $H := X(X^T X)^{-1} X^T$

$$\hat{\varepsilon}^T \hat{\varepsilon} = [(I - H)Y]^T (I - H)Y = Y^T (I - H)^T (I - H)Y = Y^T (I - H)Y$$

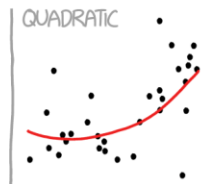
sum of squared errors (SSE)



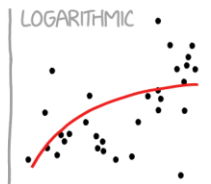
CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



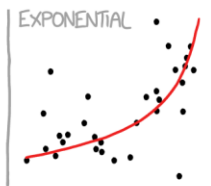
"HEY, I DID A
REGRESSION."



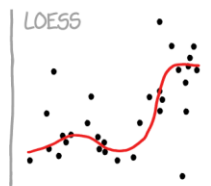
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



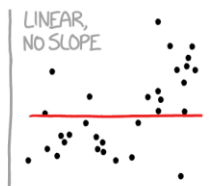
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



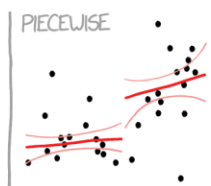
"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



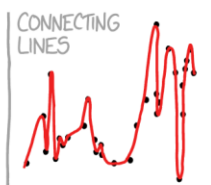
"I NEED TO CONNECT THESE
TWO LINES, BUT MY FIRST IDEA
DIDN'T HAVE ENOUGH MATH."



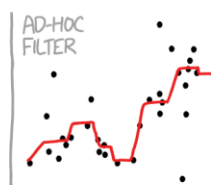
"LISTEN, SCIENCE IS HARD.
BUT I'M A SERIOUS
PERSON DOING MY BEST."



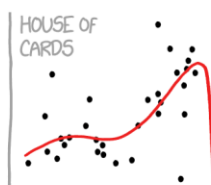
"I HAVE A THEORY,
AND THIS IS THE ONLY
DATA I COULD FIND."



"I CLICKED 'SMOOTH
LINES' IN EXCEL."



"I HAD AN IDEA FOR HOW
TO CLEAN UP THE DATA.
WHAT DO YOU THINK?"



"AS YOU CAN SEE, THIS
MODEL SMOOTHLY FITS THE-
WAIT NO NO DON'T
EXTEND IT AAAAAA!!"

Outline

Maximum Likelihood Estimation (MLE)

- idea
- examples

Linear Models

- classical model
- regularization
- extensions

Regression



$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{N} \|Y - X\beta\|^2 \right\}$$

classical model
regularizations
extensions

λ Lagrangian Multiplier

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{N} \|Y - X\beta\|^2 + \lambda \|\beta\|^1 \right\}$$

the Loss Function
 $L(X, Y, \lambda)$

L1 or **Least absolute shrinkage and selection operator**
- encourages **sparsity** of β
- reduces **overfitting**

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{N} \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \right\}$$

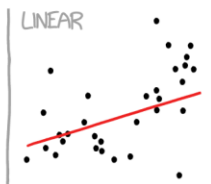
L2 or **Ridge**
- **penalizes large β**

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{N} \|Y - X\beta\|^2 + \lambda \max(0, -\beta) \right\} \quad - \text{penalizes negative } \beta$$

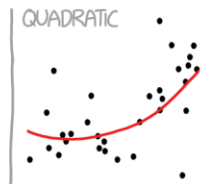
...and so on



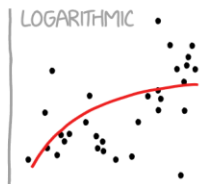
CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



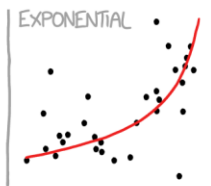
"HEY, I DID A
REGRESSION."



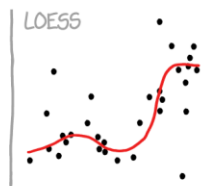
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



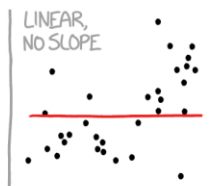
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



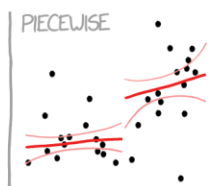
"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



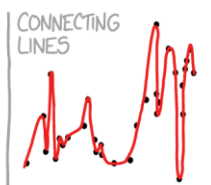
"I NEED TO CONNECT THESE
TWO LINES, BUT MY FIRST IDEA
DIDN'T HAVE ENOUGH MATH."



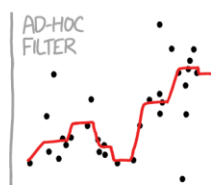
"LISTEN, SCIENCE IS HARD.
BUT I'M A SERIOUS
PERSON DOING MY BEST."



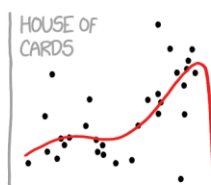
"I HAVE A THEORY,
AND THIS IS THE ONLY
DATA I COULD FIND."



"I CLICKED 'SMOOTH
LINES' IN EXCEL."



"I HAD AN IDEA FOR HOW
TO CLEAN UP THE DATA.
WHAT DO YOU THINK?"



"AS YOU CAN SEE, THIS
MODEL SMOOTHLY FITS THE-
WAIT NO NO DON'T
EXTEND IT AAAAAA!!"

Outline

Maximum Likelihood Estimation (MLE)

- idea
- examples

Linear Models

- classical model
- regularization
- extensions

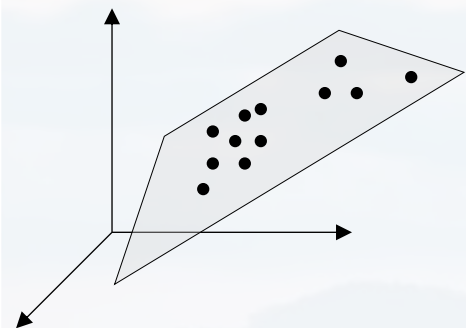
Regression



$$y_k = \beta_0 + \sum_{n=1}^N \beta_n x_n + \epsilon_k$$

the classical model is very limiting

classical model
regularizations
extensions



general: linear refers to the **factors**

$$y_k = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

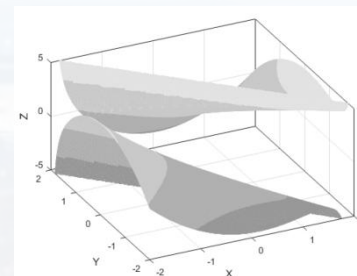
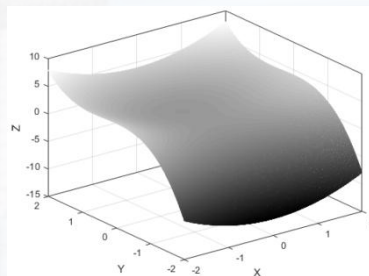
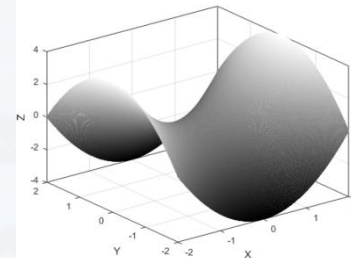
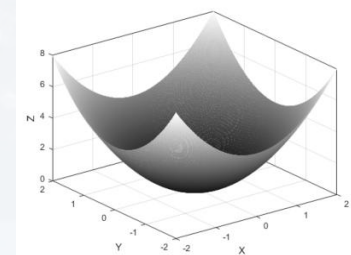
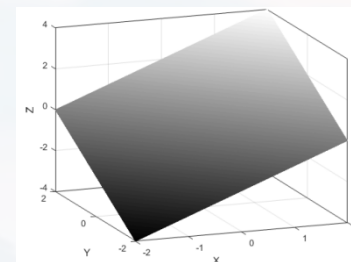
2D plane in 3D space

$$y_k = \beta_0 + \beta_1 x_1^2 + \beta_2 x_2^2$$

2D parabolic

$$y_k = \beta_0 + \beta_1 x_1^2 - \beta_2 x_2^2$$

2D hyperbolic

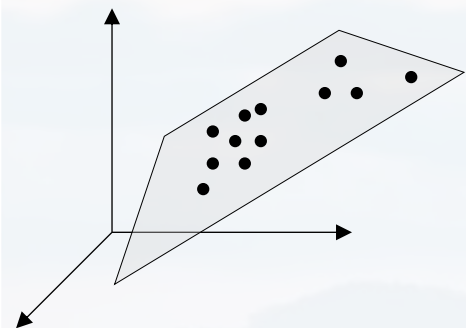


...and many more...



$$y_k = \beta_0 + \sum_{n=1}^N \beta_n x_n + \epsilon_k \quad \text{the classical model is very limiting}$$

classical model
regularizations
extensions



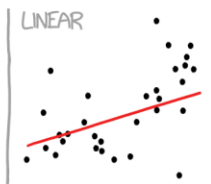
general: linear refers to the **factors**

- x_n can be subject to any **basis function** $\varphi_n(x_n)$
- the matrix X is then called **design matrix** $\phi_{i,j} = \varphi_j(x_i)$
- $\hat{\beta} = (X^T X)^{-1} X^T Y \rightarrow (\phi^T \phi)^{-1} \phi^T Y$ is called **normal equations**
- most textbooks: rows = number of observations/samples (here K)
columns = number of features (here N)
 \rightarrow one φ_j per column j
- 1D curve fit \rightarrow **one feature** x :
$$y_k = \beta_0 + \sum_{n=1}^N \beta_n \varphi_n(x_n) + \epsilon_k$$

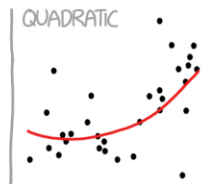
polynomial fit
$$y_k = \beta_0 + \sum_{n=1}^N \beta_n x_{k,n}^n + \epsilon_k$$



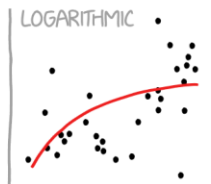
CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



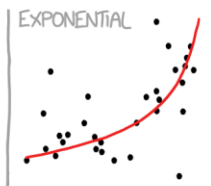
"HEY, I DID A
REGRESSION."



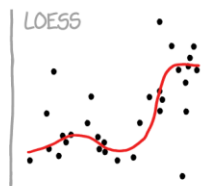
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



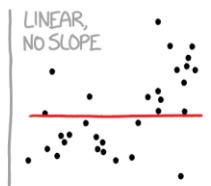
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



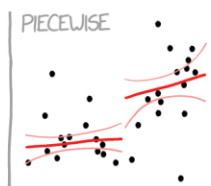
"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



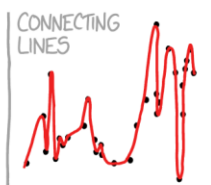
"I NEED TO CONNECT THESE
TWO LINES, BUT MY FIRST IDEA
DIDN'T HAVE ENOUGH MATH."



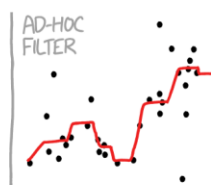
"LISTEN, SCIENCE IS HARD.
BUT I'M A SERIOUS
PERSON DOING MY BEST."



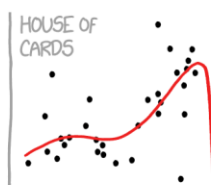
"I HAVE A THEORY,
AND THIS IS THE ONLY
DATA I COULD FIND."



"I CLICKED 'SMOOTH
LINES' IN EXCEL."



"I HAD AN IDEA FOR HOW
TO CLEAN UP THE DATA.
WHAT DO YOU THINK?"



"AS YOU CAN SEE, THIS
MODEL SMOOTHLY FITS THE-
WAIT NO NO DON'T
EXTEND IT AAAAAA!!"

Outline

Maximum Likelihood Estimation (MLE)

- idea
- examples

Linear Models

- classical model
- regularization
- extensions

Regression



We now want to apply MLE to our problem and see how it leads to the same structure for $\hat{\beta}$

$$Y = X\beta + \varepsilon$$

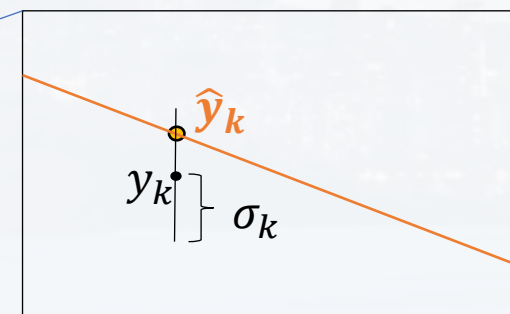
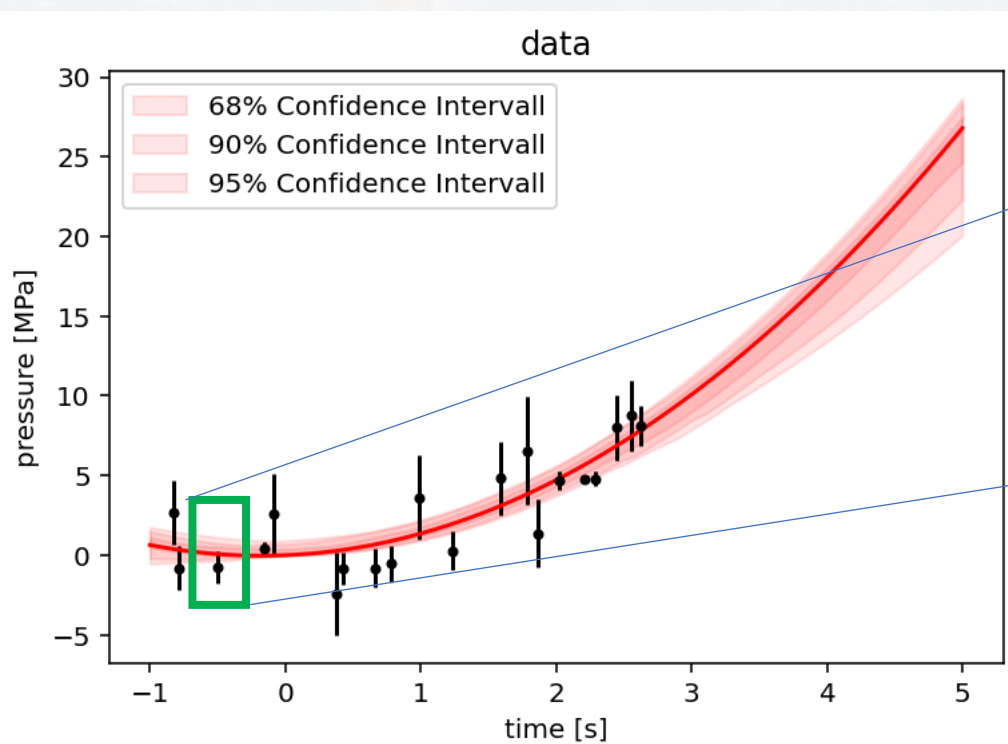
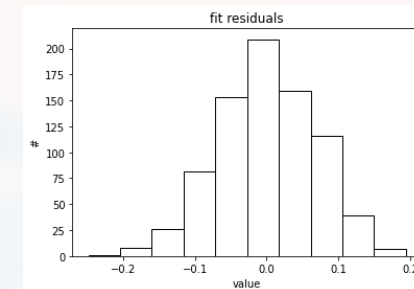
ε : additive gaussian noise around $\mu = 0$

function $y(x, \beta)$

we aim to predict a target value \hat{y}

$$\hat{y} = y(x, \beta) + \varepsilon$$

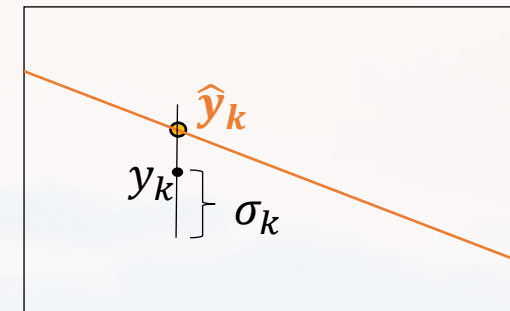
for each data point k





$$p(\hat{y}|\beta, \sigma) = \prod_{k=1}^K N(\hat{y}_k | \beta^T \phi(x_k), \sigma)$$

note: we assume $\sigma_k = \text{const} = \sigma$
for simplification



$$\ln[p(\hat{y}|\beta, \sigma)] = -\frac{K}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{k=1}^K [\hat{y}_k - \beta^T \phi(x_k)]^2$$

We want to find $\hat{\beta}$ via MLE: gradient of $\ln[p(\hat{y}|\beta, \sigma)]$ wrt β

$$\text{grad}\{\ln[p(\hat{y}|\beta, \sigma)]\}_{\beta} = \frac{1}{\sigma^2} \sum_{k=1}^K [\hat{y}_k - \beta^T \phi(x_k)] \phi(x_k)^T$$

setting the gradient to zero:

$$\hat{\beta} = (\phi^T \phi)^{-1} \phi^T \hat{y}$$

as before

for different σ_k we find

$$\chi^2 = \sum_{k=1}^K \left(\frac{\hat{y}_k - \beta^T \phi(x_k)}{\sigma_i} \right)^2$$



Thank you very much for your attention!

