Statistical modeling: the two cultures

Leo Breiman

Statistical Science, 2001

There are two cultures in the use of statistical modeling to reach conclusions from data.

- One assumes that the data are generated by a given data model.
- The other uses algorithmic models and treats the data mechanism as unknown.

The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems.

In the mid-1980s two powerful new algorithms for fitting data became available: neural nets and decision trees. A new research community using these tools sprang up. Their goal was predictive accuracy. The community consisted of young computer scientists, physicists and engineers plus a few aging statisticians. They began using the new tools in working on complex prediction problems where it was obvious that data models were not applicable: speech recognition, image recognition, nonlinear time series prediction, handwriting recognition, prediction in financial markets.

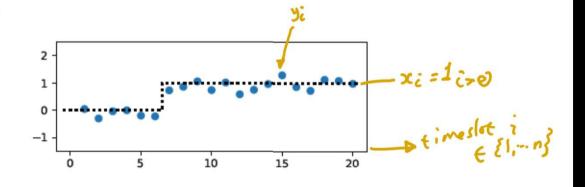
TODAY Mock exam question 1

Printouts are on the front bench, both sides.

A 0/1 signal is being transmitted. The transmitted signal at timeslot $i \in \{1, \ldots, n\}$ is $x_i \in \{0, 1\}$, and we have been told that this signal starts at 0 and then flips to 1, i.e. there is a parameter $\theta \in \{1, \ldots, n-1\}$ such that $x_i = 1_{i>\theta}$. The value of this parameter is unknown. The channel is noisy, and the received signal in timeslot i is

$$Y_i \sim x_i + \text{Normal}(0, \varepsilon^2)$$

where ε is known.



- (i) Given received signals (y_1, \ldots, y_n) , find an expression for the log likelihood, $\log \Pr(y_1, \ldots, y_n; \theta)$. Explain your working. [5 marks]
- (ii) Give pseudocode for finding the maximum likelihood estimator $\hat{\theta}$.

 Skim read for keywords. What's the topic?

log (ikelihood)
mle
pseudocode for mle

- 2. Look for question words. What is it asking you to do?

 Find expression for me

 give pseudocade for fix.
 - 3. Think through the course. What sections are relevant?

Sc1. specifying sectly models

5 marks

Indicator functions

ex 1.3.6, §2.2.1

Indicator functions are a notational trick that lets us turn if/then conditions into algebra. The indicator function is defined to be

$$1_A = \begin{cases} 1 & \text{if statement } A \text{ is true} \\ 0 & \text{if statement } A \text{ is false} \end{cases}$$

Observe that they turn *and* into *multiplication*:

$$1_{A \text{ and } B} = 1_A \cdot 1_B$$

In IA Discrete Maths [lecture 17] you used the indicator for a subset:

$$\chi_S(a) = 1_{a \in S}$$

INS
$$Pr(y_1,...,y_n; \theta)$$

$$= \sum_{i=1}^{n} log Pr(y_i; \theta) \quad modelling the Y_i as independent.$$

$$= \sum_{i=1}^{n} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; \theta)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad Pr_{y_i}(y_i; x_i)$$

$$= \sum_{i} log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - x_i)^2/2\epsilon^2} \right\} \quad$$

• be a constant $-\xi, n$ • be a dumny variable -i $= \sum_{i=1}^{n} \log \left\{ \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-(y_i^2 - 1_{i,n})^2/2\epsilon^2} \right\}$ $= -\frac{n}{2} \log (2\pi\epsilon^2) - \frac{1}{2\epsilon^2} \sum_{i=1}^{n} (y_i^2 - 1_{i,n})^2.$

(ii) We need to find max of the log likelihood we just derived.

def loglik (0): retuin...

scipy-optimize. frain (lambda 0: -loglik (0) init-quess)

slipy-optimize. frain (lambda 0: -loglik (0) init-quess)

slipy-optimize. frain (lambda 0: -loglik (0) init-quess)

vols = [(loglik(0), 0)] for $0 \in \{1, ---, n-1\}$] $(-, \vec{\theta}) = max(vols)$

- (b) I have been monitoring average annual river levels for many years, and I have collected a dataset (z_0, \ldots, z_n) where z_i is the level in year i since I started monitoring. I believe that for the first few years the level each year was roughly what it was the previous year, plus or minus some random variation; but that some year a drought started, and since then the level has decreased on average each year. I would like to estimate when the drought started. I do not know the other parameters.
 - (i) Propose a probability model for my dataset.
 - (ii) Explain how to fit your model.

This "propose a probability model" is open-ended and scary. How should we even begin to think about it?

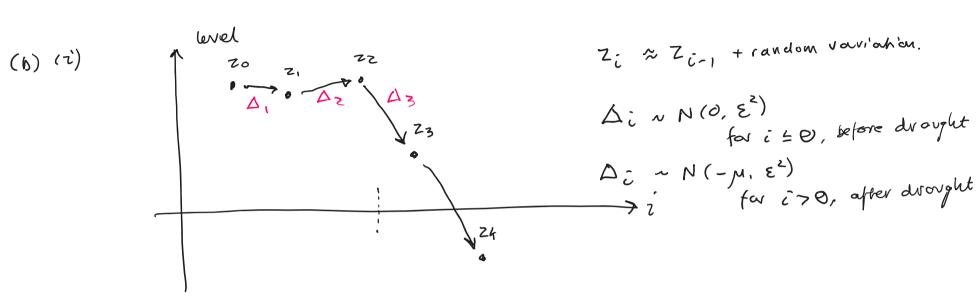
- Skim read for keywords. What's the topic?
- 2. Look for question words. What is it asking you to do?

[5 marks]

5 marks

- 3. Think through the course. What sections are relevant?
- 4. Read the whole question. What's the link?

→ Part (a) gave us a hammer. Can we see part (b) as a naîl?



Model: $\Delta_i \sim N(-\mu \mathbf{1}_{i70}, \xi^2)$ where $\Delta_i = 2i - 2i - 1$ μ, θ, ξ unknown.

(ii)
$$\log \Pr(\delta_1, ..., \delta_n; M, \theta, \epsilon)$$

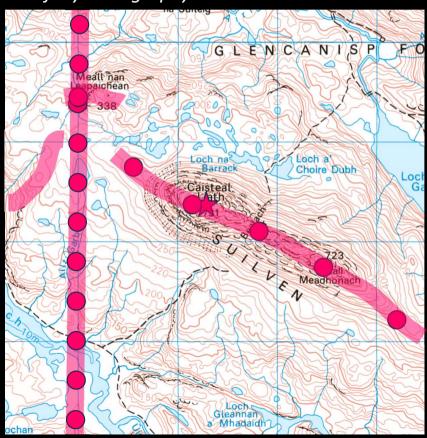
= $-\frac{n}{2} \log (2\pi\epsilon^2) - \frac{1}{2\epsilon^2} \sum_{i=1}^{n} (\delta_i + M \mathbf{1}_{i>0})^2$

we want to estimate o.

THEOREM

$$\max_{x,y} f(x,y) = \max_{x} \left\{ \max_{y} f(x,y) \right\}$$

Proof: By cartography.



 $\max_{x,y} f(x,y)$ scans the entire area looking for the highest point

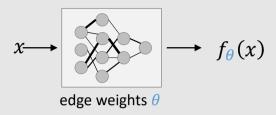
 $\max_{y} f(x, y)$ scans a line of fixed x, looking for the highest point

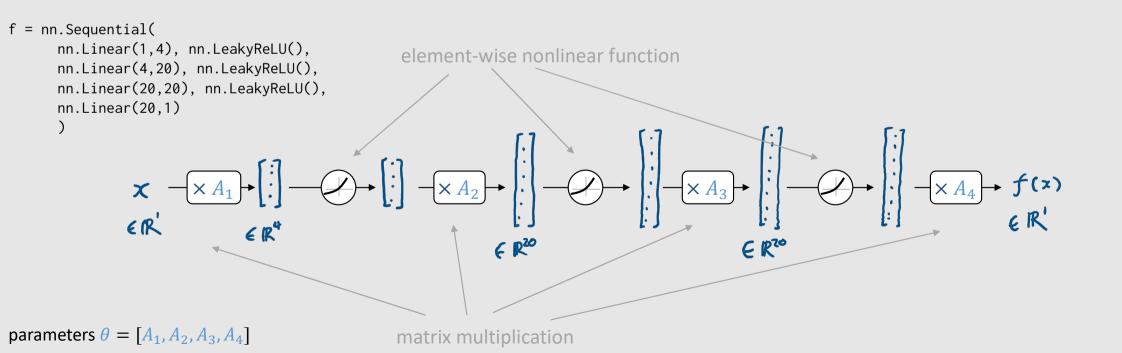
$$\max_{x} \left\{ \max_{y} f(x, y) \right\}$$

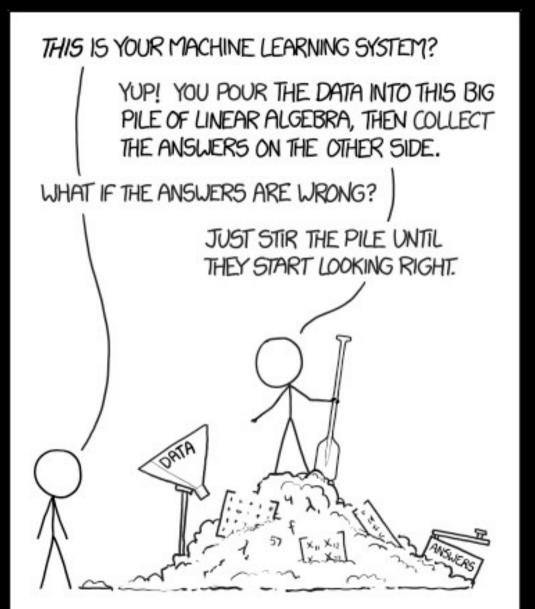
look at all the vertical slices, and for each find the highest point then pick the highest out of all these Neural networks are probability models*

* non-examinable

A neural network is a deterministic parametric function, typically consisting of matrix multiplication and nonlinear functions.







https://xkcd.com/1838 CC BY-NC 2.5

PREDICTION MINDSET IN MACHINE LEARNING

We're given a labelled dataset.

We want to learn to predict the label.

We do this by minimizing a loss function.

ground-truth labels: Let y_i be the actual observed temperature at time t_i



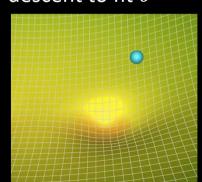
it's our job as modellers to find θ so as to minimize prediction error, e.g. pick θ to minimize

$$\sum_{i} L(y_i, \hat{y}_{\theta}(t_i))$$

where

$$L(obs, pred) = (obs - pred)^2$$

We can use gradient descent to fit θ



All that matters about the neural network is

- it's a parametric function $t \mapsto \hat{y}_{\theta}(t)$
- it's differentiable with respect to its parameters θ

The loss function L is arbitrary, but it had better be differentiable with respect to pred.

This course teaches a different way to think of modelling ...

PREDICTION MINDSET



Goal:

minimize prediction error

PROBABILITY MODELLER'S VIEW

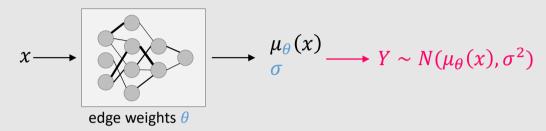


Goal:

maximize the log likelihood of the dataset

Q. How does a neural network output a random variable?

Easy mode: have it output the parameters of a distribution



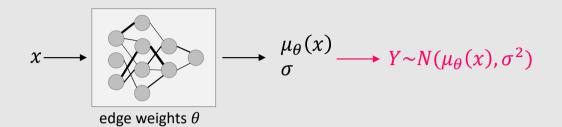
❖ Hard mode: give it some noise as input

$$Z \sim N(0,1)$$

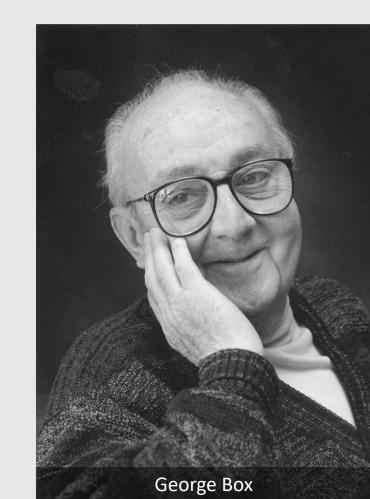
$$X \longrightarrow Y = f_{\theta}(x, Z)$$
edge weights θ

(We then have to compute the likelihood $Pr_Y(y)$. This is what makes these easy or hard.)

Q. Why does the response have to be Gaussian?



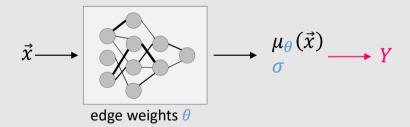
- ❖ It doesn't!
- All models are wrong
- ... but some are useful



Q. How would we specify a probability model that uses several predictor variables?

$$Y \sim N(\mu_{\theta}(\vec{x}), \sigma^2)$$

where \vec{x} is a tuple of predictor variables, $\vec{x} = (x_1, ..., x_d)$, and $\mu_{\theta} : \mathbb{R}^d \to \mathbb{R}$ is a differentiable function



- \bullet Easy mode: when $\vec{x} \in \mathbb{R}^d$ for fixed dimension d
- \clubsuit Hard mode: when \vec{x} is an arbitrary-length sequence, e.g. a text prompt

In this temperature-forecasting example, there's not very much difference in practice between the prediction mindset and the probability modelling view.

In other examples, the prediction mindset runs into problems ...

Neural network classification

The MNIST database of handwritten images consists of records (x_i, y_i) where $x_i \in \mathbb{R}^{28 \times 28}$ is a greyscale image with 28×28 pixels, and $y_i \in \{0, ..., 9\}$ is the digit.

We'd like to predict the digit, given an image. How might we learn to do this?

Data from http://yann.lecun.com/exdb/mnist/



 \int predicted
label $\hat{y}_{\theta}(x)$

ground truth:

Let y_i be the actual observed label in the dataset

it's our job as modellers to find $\boldsymbol{\theta}$ so as to maximize prediction accuracy, i.e.

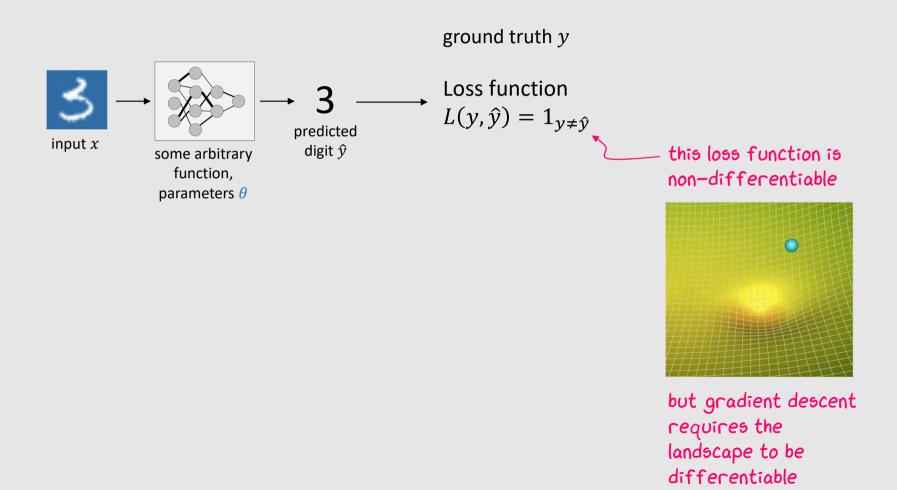
pick θ to minimize

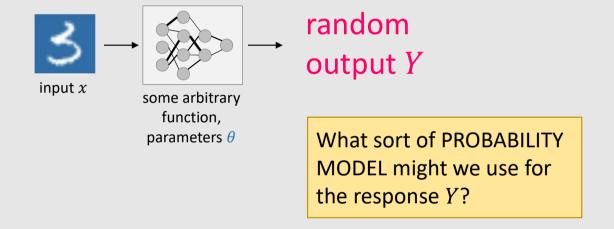
$$\sum_{i} L(y_i, \hat{y}_{\theta}(x_i))$$

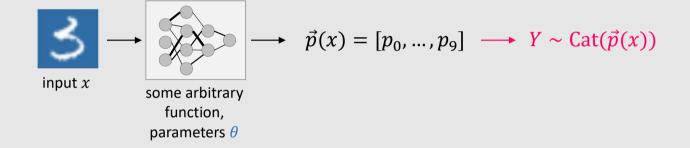
where

$$L(obs, pred) = 1_{obs \neq pred}$$

Q. How would we train this neural network with gradient descent?





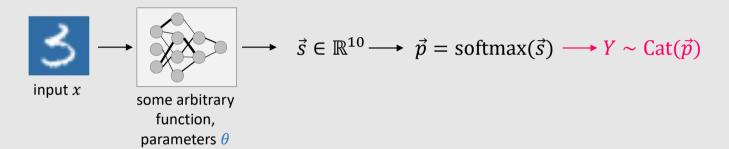


How can we make sure that \vec{p} is a valid probability vector?

(We need $p_i \in [0,1]$ for each i, and $\Sigma_i p_i = 1$.)

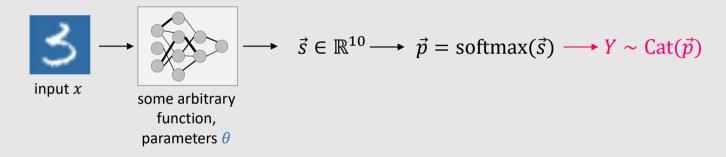
Softmax function:

$$p_k = \frac{e^{s_k}}{\sum_{\ell=0}^9 e^{s_\ell}}$$



How do we fit this model to the data?

0. Probability model for a single datapoint:



1. Likelihood of a ground-truth datapoint *y*:

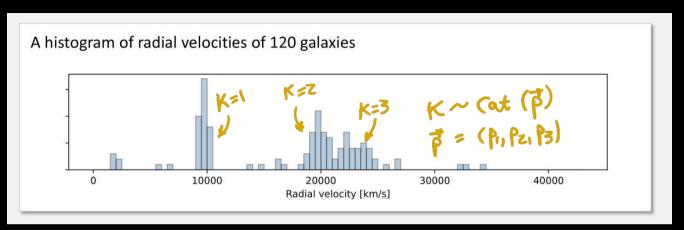
$$Pr_{V}(y; x, \theta) =$$

2. Log likelihood of the whole dataset:

$$\log \Pr(y_1, ..., y_n) = \ldots$$
 we end up with the famous "softmax cross-entropy loss function"

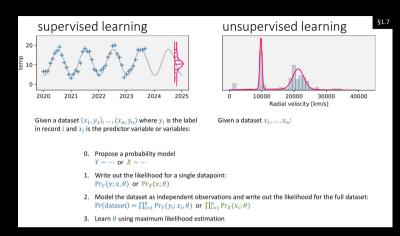
3. Find the maximum likelihood estimator $\hat{\theta}$ using gradient descent

Categorical random variable (lec3)



Softmax reparameterization trick (lec2)

Supervised learning (lec3)



PREDICTION MINDSET



Goal:

pick θ to maximize prediction accuracy

but we have to fudge the loss function to make it differentiable 🤪



PROBABILITY MODELLER'S VIEW



Goal:

fit my Y_{θ} model to the dataset by maximizing log likelihood

Statistical modeling: the two cultures Leo Breiman

Statistical Science, 2001

There are two cultures in the use of statistical modeling to react anclusions from data.

- One assu. s that the data are generated by a given date...del
- The other uses porithmic models and treats the demechanism as unknown.

The statistical community as been committed to it almost exclusive use of data models. This commitment has to irrelevant cory, questionable conclusions, and has kept statisticians from working on the regular geof interesting current problems.

In the mid-1980s two powerful new ms for fitting data became available: neural nets and decision trees community using these tools sprang .ew research up. Their goal was prediction munity consisted of young accuracy. The co cists and engineers plu computer scientists, ph few aging statisticians. They ools in working on complex procession problems where it began using the ne ata models were not applicable: speed was obvious the ecognition, image recognition conlinear time series prediction, handwriting remition, prediction in finar a markets.

Neural networks can be viewed as probability models for data.

- This view solves problems with the "prediction mindset"
- It unifies prediction, description e.g. clustering, and generative AI
- It opens the door to deeper thinking about generalizability

■ Office hours in the cafe area 1–1.30pm today

