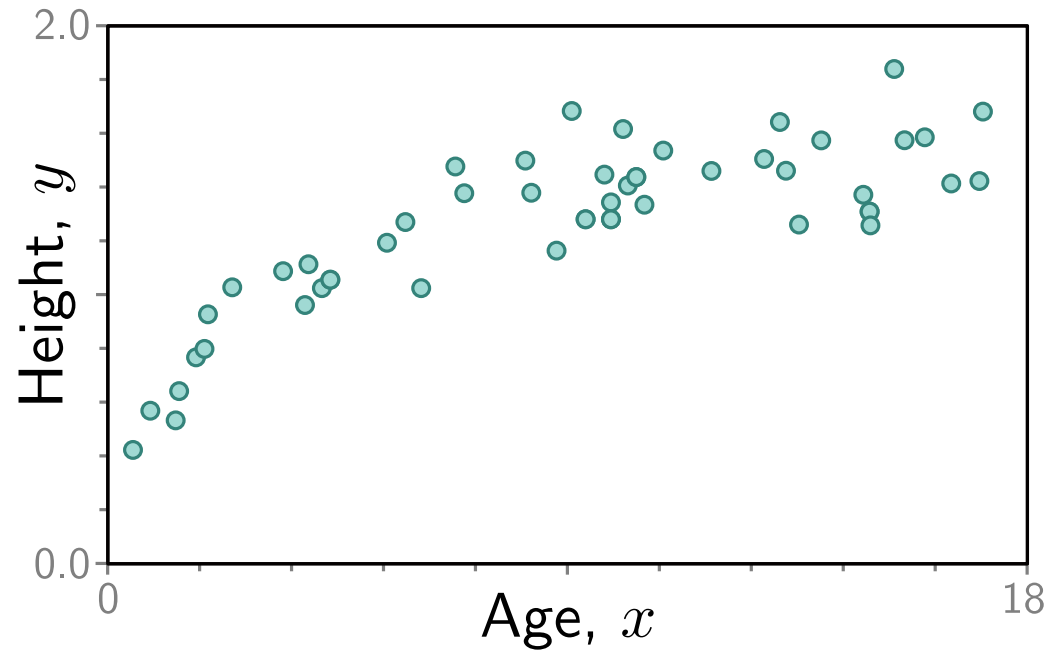# Lecture 05
# Loss Functions
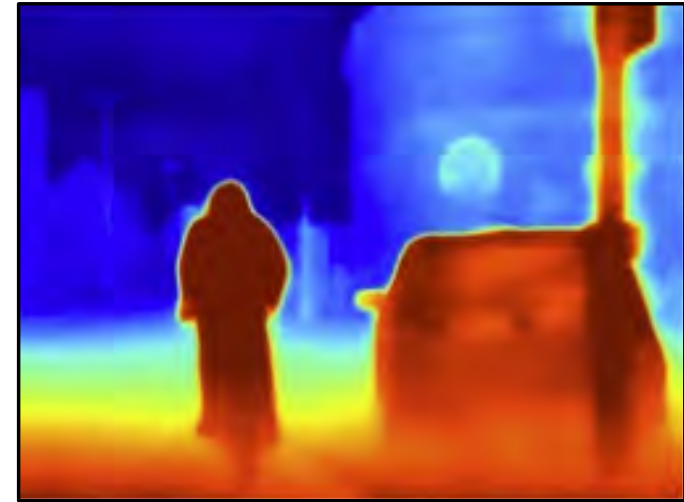# (and probability models)

DL4DS – Spring 2024

# Recap

- So far, we talked about *linear regression*, *shallow neural networks* and *deep neural networks*

- Each have parameters, $\phi$, that we want to choose for a *best possible mapping between input and output* training data

- A *loss function* or *cost function*, $L[\phi]$, returns a single number that describes a mismatch between $f[x_i, \phi]$ and the ground truth outputs, $y_i$.
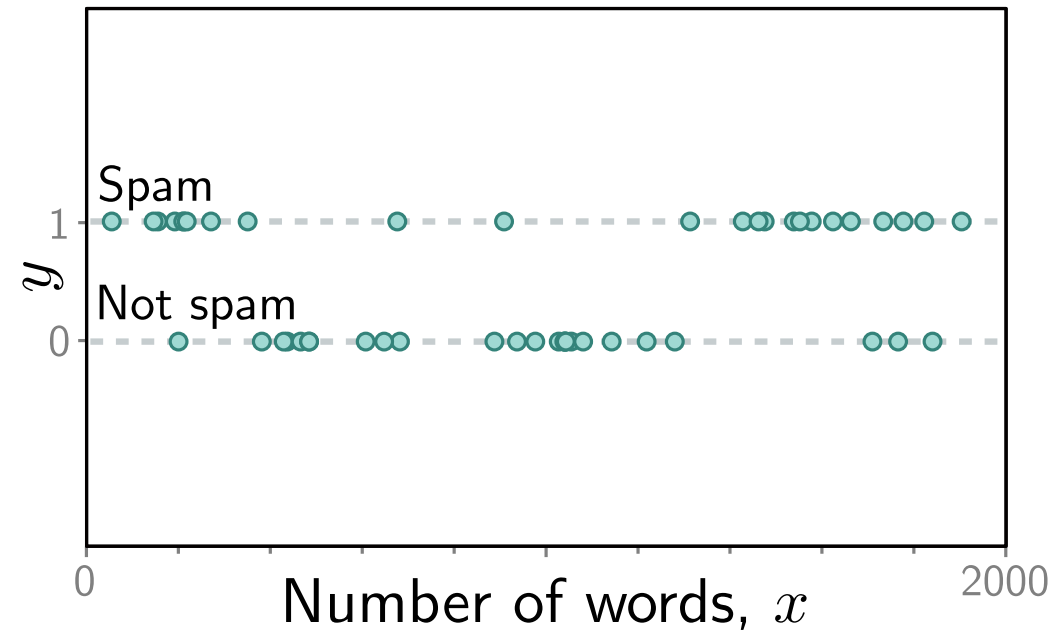
We need to find a loss function that works with…

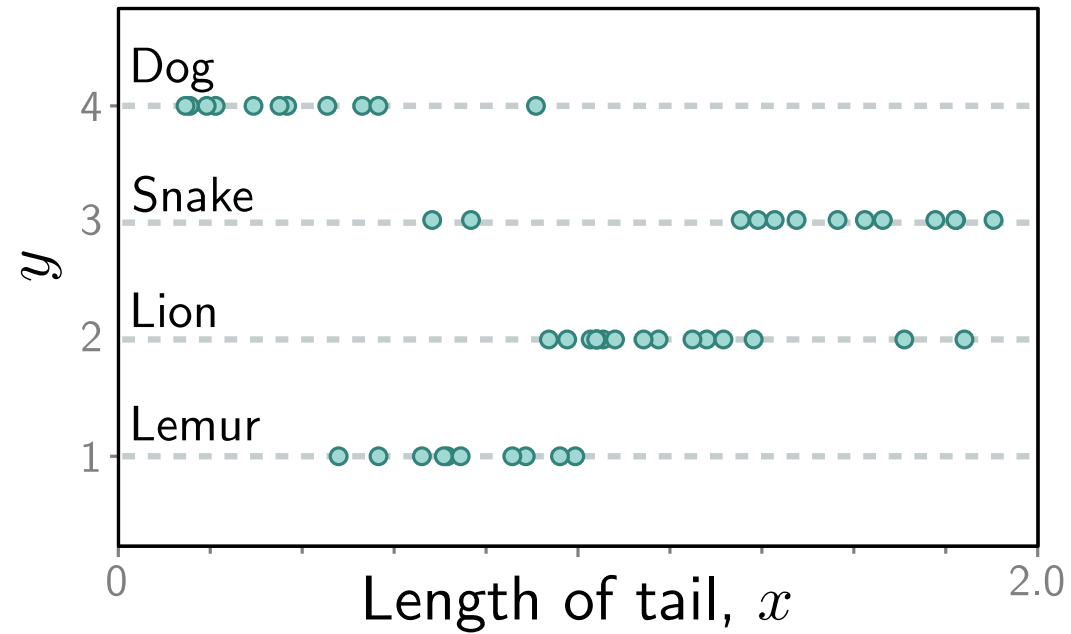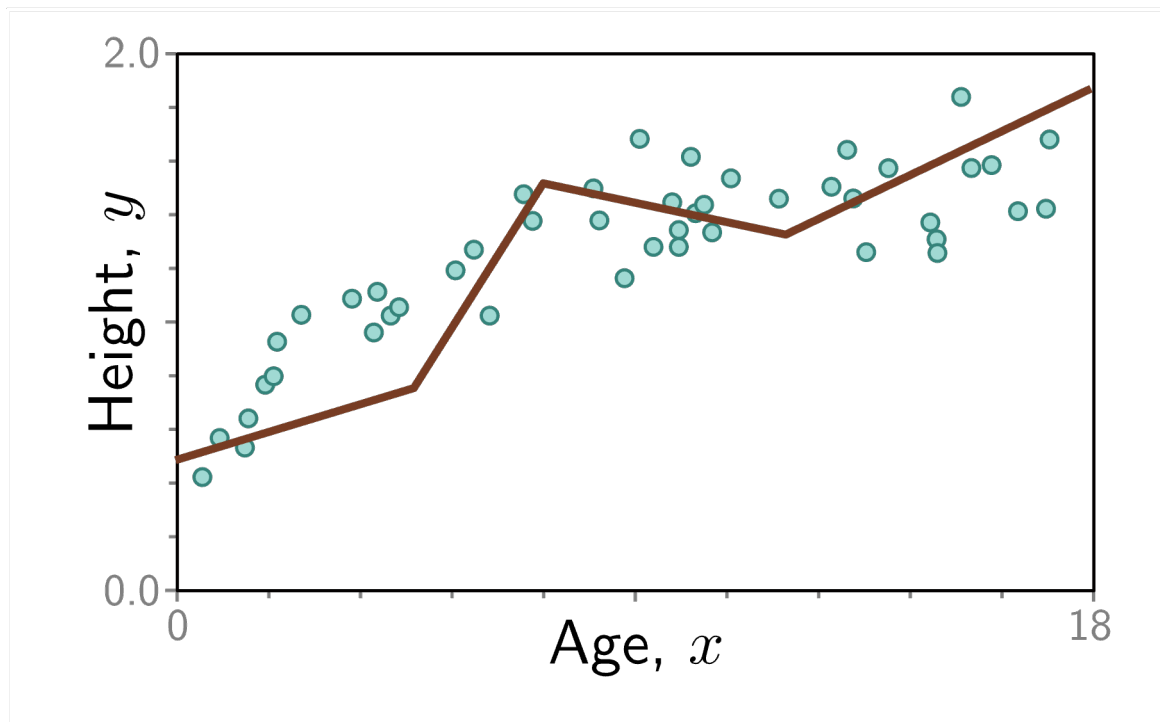# Univariate and Multivariate Regression




Depth Map

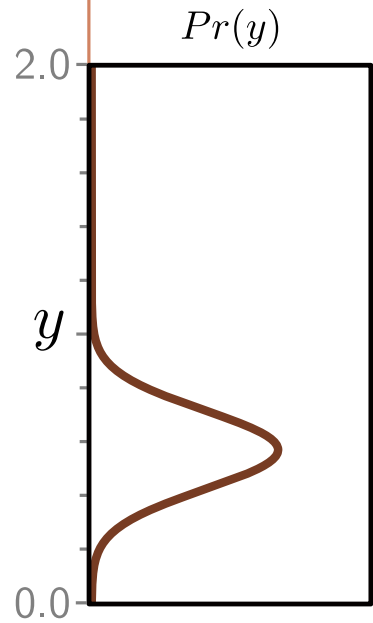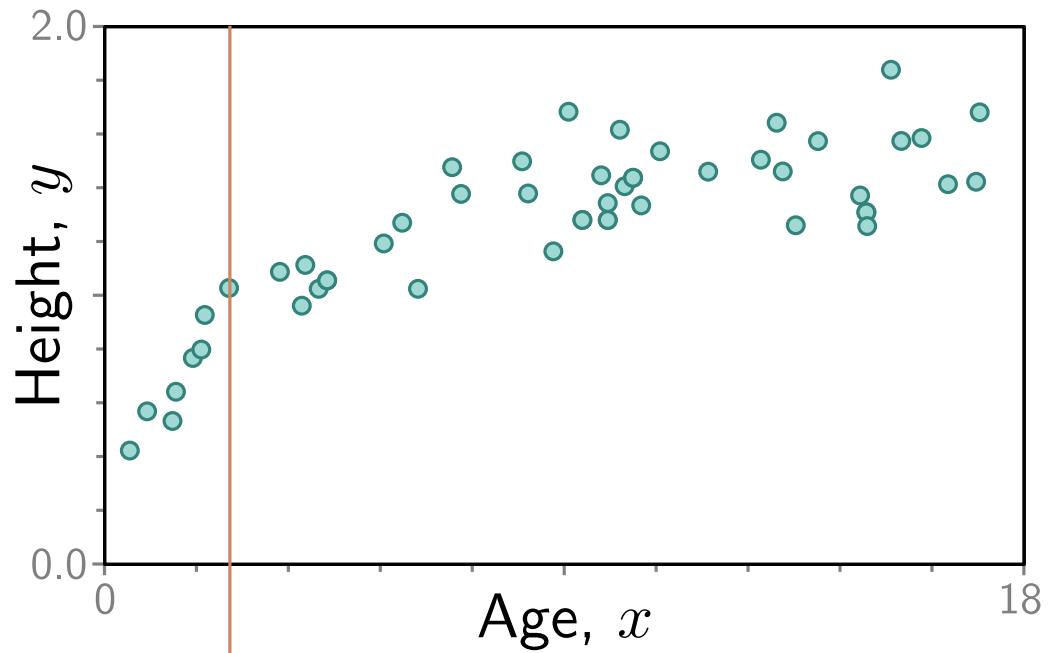# Binary Classification

# Multiclass Classification

So far, we thought about fitting a model to the data...

Alternatively, we can think about fitting a *probability model* to the data.

$$\Pr(y|x)$$

Why?

Alternatively, we can think about fitting a *probability model* to the data.

$$\text{Pr}(y|x)$$

Why?

Because this provides a *framework* to build loss functions for other prediction types...

Alternatively, we can think about fitting a *probability model* to the data.

$$\mathrm{Pr}(y|x)$$

Why?

Because this provides a *framework* to build loss functions for other prediction types...

... and justifies least squares for real-valued regression models.

# Brief Probability Review

- Random variables, e.g. $x$ and $y$
- $\Pr(x)$ is a probability distribution over $x$
- $0 \leq \Pr(x) \leq 1$
- $\int_x \Pr(x)\, dx = 1$  or  $\sum_i \Pr(x_i) = 1$
- $\Pr(x, y) = \Pr(x) \cdot \Pr(y)$ when $x$ and $y$ are independent
- $\Pr(x \mid y) \Pr(y) = \Pr(x, y) = \Pr(y \mid x) \Pr(x)$
- And…

# Joint and Marginal Probability Distributions

a)

Joint Distribution

$Pr(x, y)$

c)

$Pr(y)$

Marginal distribution

$$\Pr(y) = \int_x \Pr(x, y)\, dx$$

b)

Marginal distribution

$$\Pr(x) = \int_y \Pr(x, y)\, dy$$

# Conditional Probabilities

a)

$Pr(x, y)$



b)

$Pr(x|y{=}3.0)$



$$\int_x \Pr(x \mid y = 3.0)\, dx = 1$$

c)

$Pr(x|y{=}{-}1.0)$



$$\int_x \Pr(x \mid y = -1.0)\, dx = 1$$

Continuous
Pr($y|x$)

Continuous
Pr($y|x$)

Continuous
Pr(*y*|*x*)

Continuous
Pr($y|x$)

Continuous
Pr($y|x$)

Spam

$y$

Not spam

Number of words, $x$

0

2000

1

0

Spam

Not spam

$y$

Number of words, $x$

0　　　　　　　　　　　　　　　　2000

1

0

$Pr(y)$

Discrete
$\Pr(y|x)$

$y$

1

0

20

Spam

$y$

1

Not spam

0

Number of words, $x$

0                                    2000

$Pr(y)$

Discrete
$Pr(y|x)$

$y$

1

0

Spam

Not spam

$y$

Number of words, $x$

0                                                    2000

$Pr(y)$

$y$

Discrete
Pr($y|x$)

Spam

Not spam

$y$

Number of words, $x$

0

2000

1

0

$Pr(y)$

$y$

1

0

Discrete
Pr($y$|$x$)

Spam

$y$

1

Not spam

0

Number of words, $x$

0                    2000

$Pr(y)$

$y$

1

0

0          1

Discrete
$Pr(y|x)$

Dog

Snake

Lion

Lemur

$y$

Length of tail, $x$

0          2.0

$Pr(y)$

Discrete
Pr($y|x$)

$y$

4

3

2

1

Dog

Snake

Lion

Lemur

$y$

Length of tail, $x$

$Pr(y)$

Discrete
$Pr(y|x)$

$y$

Dog

Snake

Lion

Lemur

$y$

Length of tail, $x$

0

2.0

$Pr(y)$

$y$

Discrete
Pr($y|x$)

28

Dog

Snake

Lion

Lemur

$y$

Length of tail, $x$

$Pr(y)$

$y$

Discrete
Pr($y|x$)

a)

$Pr(y|x=2)$      $Pr(y|x=7)$

a)

b)

$Pr(y|x=2)$   $Pr(y|x=7)$

$Pr(y|x=2)$   $Pr(y|x=7)$

a)

b)

c)

$Pr(y|x=2)$   $Pr(y|x=7)$

$Pr(y|x=2)$   $Pr(y|x=7)$

$Pr(y|x=2)$   $Pr(y|x=7)$

# Loss function

- Training dataset of $I$ pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$$

- Loss function or cost function measures how bad model is:

$$L\left[\boldsymbol{\phi}, \underbrace{\mathrm{f}[\mathbf{x}, \boldsymbol{\phi}]}_{\text{model}}, \underbrace{\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}}_{\text{train data}}\right]$$

# Loss function

- Training dataset of $I$ pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$$

- Loss function or cost function measures how bad model is:

or for short:

$$L[\boldsymbol{\phi}]$$ ← Returns a scalar that is smaller when model maps inputs to outputs better

# Training

- Loss function:

$$L\left[\phi\right]$$

Returns a scalar that is smaller when model maps inputs to outputs better

- Find the parameters that minimize the loss:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}}\left[\mathrm{L}\left[\phi\right]\right]$$

# Example: 1D Linear regression loss function

Loss, $L = 7.11$

$\phi_0 = 0.4, \phi_1 = 0.2$

Output, $y$

Input, $x$

Loss function:

$$L[\phi] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \phi] - y_i)^2$$

$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression training



This technique is known as gradient descent

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1:  univariate regression
- Example 2:  binary classification
- Example 3:  multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Maximum Likelihood Estimation

- In statistics, *maximum likelihood estimation (MLE)* is a method of *estimating the parameters* of an *assumed probability distribution*, *given some observed data*.

- This is achieved by *maximizing a likelihood function* so that, under the assumed statistical model, *the observed data is most probable*.

# How do we do this?

- Model predicts output *y* given input *x*

# How do we do this?

- ~~Model predicts output *y* given input *x*~~

# How do we do this?

- ~~Model predicts output *y* given input *x*~~
- Model predicts a conditional probability distribution:

$$Pr(\mathbf{y}|\mathbf{x})$$

over outputs **y** given inputs **x**.
- Define and minimize a loss function that makes the outputs have high probability

# How can a model predict a probability distribution? ➤ Parametric Models

1. Pick a known distribution (e.g., normal distribution) to model output y with parameters $\boldsymbol{\theta}$

   e.g., the normal distribution $\boldsymbol{\theta} = \{\mu, \sigma^2\}$



2. Use model to predict parameters $\boldsymbol{\theta}$ of probability distribution

# Maximize the joint, conditional probability

- We know we picked a good model and the right parameters when the joint conditional probability is high for the observed (e.g. training) data.

$$\text{Pr}(y_1, y_2, \dots, y_I \mid x_1, x_2, \dots, x_I)$$

# Two simplifying assumptions

Identically distributed (the form of the probably distribution is the same for each input/output pair)

$$\Pr(y_1, y_2, \ldots, y_I \mid x_1, x_2, \ldots, x_I) = \prod_{i=1}^{I} \Pr(y_i \mid x_i)$$

Independent

*Independent and identically distributed (i.i.d)*

# Maximum likelihood criterion

$$\hat{\phi} = \operatorname*{argmax}_{\phi} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{x}_i) \right]$$

$$= \operatorname*{argmax}_{\phi} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \boldsymbol{\theta}_i) \right]$$

$$= \operatorname*{argmax}_{\phi} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right]$$

$\theta_i$ are the parameters of the probability distribution

$\phi$ are the parameters of the neural network, e.g.

$$\theta_i = \mathrm{f}[\mathrm{x}_i, \phi]$$

When we consider this probability as a function of the parameters $\phi$, we call it a likelihood.

# Problem:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right]$$

- The terms in this product might all be small
- The product might get so small that we *can't* easily represent it in fixed precision arithmetic

# The log function is monotonic



Maximum of the logarithm of a function is in the same place as maximum of function

# Maximum log likelihood

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmax}} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}]) \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmax}} \left[ \log \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}]) \right] \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmax}} \left[ \sum_{i=1}^{I} \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}]) \right] \right]$$

Now it's a sum of terms, so doesn't matter so much if the terms are small

# Minimizing negative log likelihood

- By convention, we minimize things (i.e., a loss)

$$\hat{\phi} = \underset{\phi}{\mathrm{argmax}} \left[ \sum_{i=1}^{I} \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]$$

$$= \underset{\phi}{\mathrm{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]$$

$$= \underset{\phi}{\mathrm{argmin}} \left[ \mathrm{L}[\phi] \right]$$

# Inference

- But now we predict a probability distribution

- We need an actual prediction (point estimate)

- Find the peak of the probability distribution (i.e., mean for normal)

$$\hat{y} = \hat{\mu} = \underset{y}{\mathrm{argmax}}[\mathrm{Pr}(y \mid \mathbf{f}[\mathbf{x}, \phi])]]$$

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1:  univariate regression
- Example 2:  binary classification
- Example 3:  multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Recipe for loss functions

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

# Recipe for loss functions

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$.

# Recipe for loss functions

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$.

3. To train the model, find the network parameters $\hat{\boldsymbol{\phi}}$ that minimize the negative log-likelihood loss function over the training dataset pairs $\{\mathbf{x}_i, \mathbf{y}_i\}$:

$$\hat{\boldsymbol{\phi}} = \underset{\phi}{\operatorname{argmin}} \left[ L[\boldsymbol{\phi}] \right] = \underset{\phi}{\operatorname{argmin}} \left[ - \sum_{i=1}^{I} \log \left[ Pr(\mathbf{y}_i|\mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}]) \right] \right]. \qquad (5.7)$$

# Recipe for loss functions

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

2. Set the machine learning model $\mathbf{f}[\mathbf{x},\boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x},\boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x},\boldsymbol{\phi}])$.

3. To train the model, find the network parameters $\hat{\boldsymbol{\phi}}$ that minimize the negative log-likelihood loss function over the training dataset pairs $\{\mathbf{x}_i, \mathbf{y}_i\}$:

$$\hat{\boldsymbol{\phi}} = \underset{\phi}{\operatorname{argmin}}\left[L[\boldsymbol{\phi}]\right] = \underset{\phi}{\operatorname{argmin}}\left[-\sum_{i=1}^{I}\log\left[Pr(\mathbf{y}_i|\mathbf{f}[\mathbf{x}_i,\boldsymbol{\phi}])\right]\right]. \qquad (5.7)$$

4. To perform inference for a new test example $\mathbf{x}$, return either the full distribution $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x},\hat{\boldsymbol{\phi}}])$ or the maximum of this distribution.

# Let's apply this recipe to

- Example 1: Real valued univariate regression

- Example 2: Binary Classification

- Example 3: Multiclass Classification

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Example 1: univariate regression

Real world input

Model input

Model

Model output

Real world output

6000 square feet,
4 bedrooms,
previously sold for
$235K in 2005,
1 parking spot.

$$\begin{bmatrix} 6000 \\ 4 \\ 235 \\ 2005 \\ 1 \end{bmatrix}$$

Supervised learning model

$[340]$

Predicted price
is $340k

# Example 1: univariate regression

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

- Predict scalar output: $y \in \mathbb{R}$

- Sensible probability distribution:
  - Normal distribution

$$Pr(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y-\mu)^2}{2\sigma^2}\right]$$

# Example 1: univariate regression

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \phi]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \phi]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \phi])$.

$$Pr(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y-\mu)^2}{2\sigma^2}\right]$$

In this case, just the mean

$$Pr(y|\mathbf{f}[\mathbf{x}, \phi], \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y-\mathbf{f}[\mathbf{x}, \phi])^2}{2\sigma^2}\right]$$

Just learn the mean, $\mu$, and assume the variance is fixed,.

# Example 1: univariate regression

3. To train the model, find the network parameters $\hat{\phi}$ that minimize the negative log-likelihood loss function over the training dataset pairs $\{\mathbf{x}_i, \mathbf{y}_i\}$:

$$L[\phi] = -\sum_{i=1}^{I} \log\left[Pr(y_i|\mathrm{f}[\mathbf{x}_i, \phi], \sigma^2)\right]$$

$$= -\sum_{i=1}^{I} \log\left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \phi])^2}{2\sigma^2}\right]\right]$$

$$\hat{\boldsymbol{\phi}} = \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$\hat{\boldsymbol{\phi}} = \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] + \log \left[ \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$\log[a \cdot b] = \log[a] + \log[b]$$

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] + \log \left[ \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

$$\hat{\boldsymbol{\phi}} = \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] + \log \left[ \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

Just a constant offset



68

$$\hat{\boldsymbol{\phi}} = \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] + \log\left[ \exp\left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

$$= \operatorname*{argmin}_{\boldsymbol{\phi}} \left[ -\sum_{i=1}^{I} -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[ -\frac{(y_i - f[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$
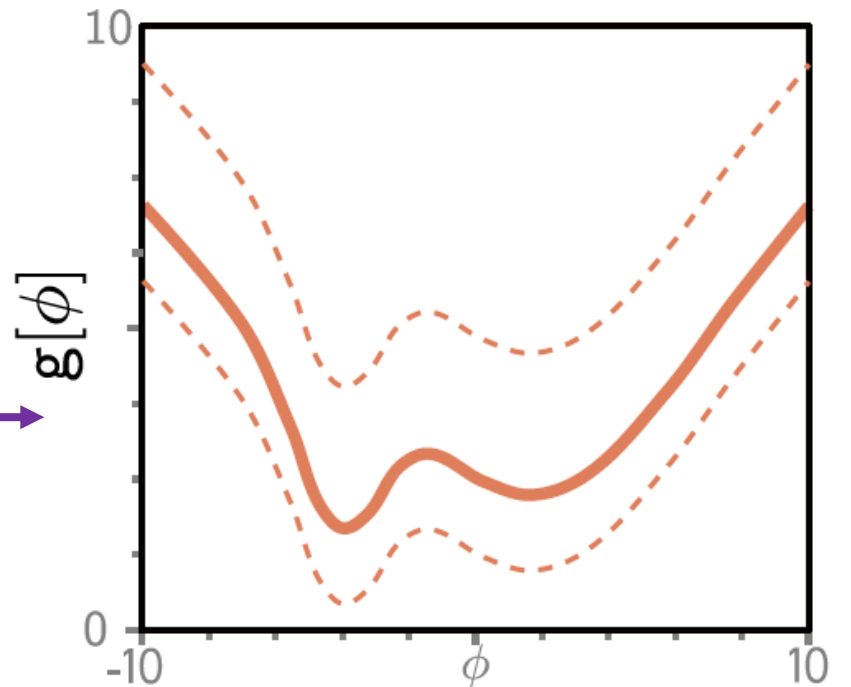
$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] + \log \left[ \exp\left[ -\frac{(y_i - f[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - f[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} -\frac{(y_i - f[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

Just dividing by a positive constant

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] + \log \left[ \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right]$$

$$= \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ \sum_{i=1}^{I} (y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2 \right], \qquad \longleftarrow \qquad \text{Least squares!}$$

# Least squares

# Negative log likelihood

a) $\sum_i (y_i - \mathsf{f}[x_i, \phi])^2 = 0.19$



c) $-\sum_i \log \left[ Pr(y_i | \mathsf{f}[x_i, \phi], \sigma^2 \right] = -6.57$

# Least squares

# Maximum likelihood



b) $\sum_i (y_i - \mathsf{f}[x_i, \phi])^2 = 10.22$

d) $-\sum_i \log \left[ Pr(y_i | \mathsf{f}[x_i, \phi], \sigma^2 \right] = 497.37$

Input, $x$

Input, $x$

# Example 1: univariate regression

4. To perform inference for a new test example $\mathbf{x}$, return either the full distribution $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \hat{\boldsymbol{\phi}}])$ or the maximum of this distribution.

Full distribution:

$$Pr(y|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}], \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y - \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])^2}{2\sigma^2}\right]$$

Max probability:

$$\hat{y} = \hat{\mu} = \mathbf{f}[\mathbf{x}|\boldsymbol{\phi}]$$

# Estimating variance

- Perhaps surprisingly, the variance term disappeared:

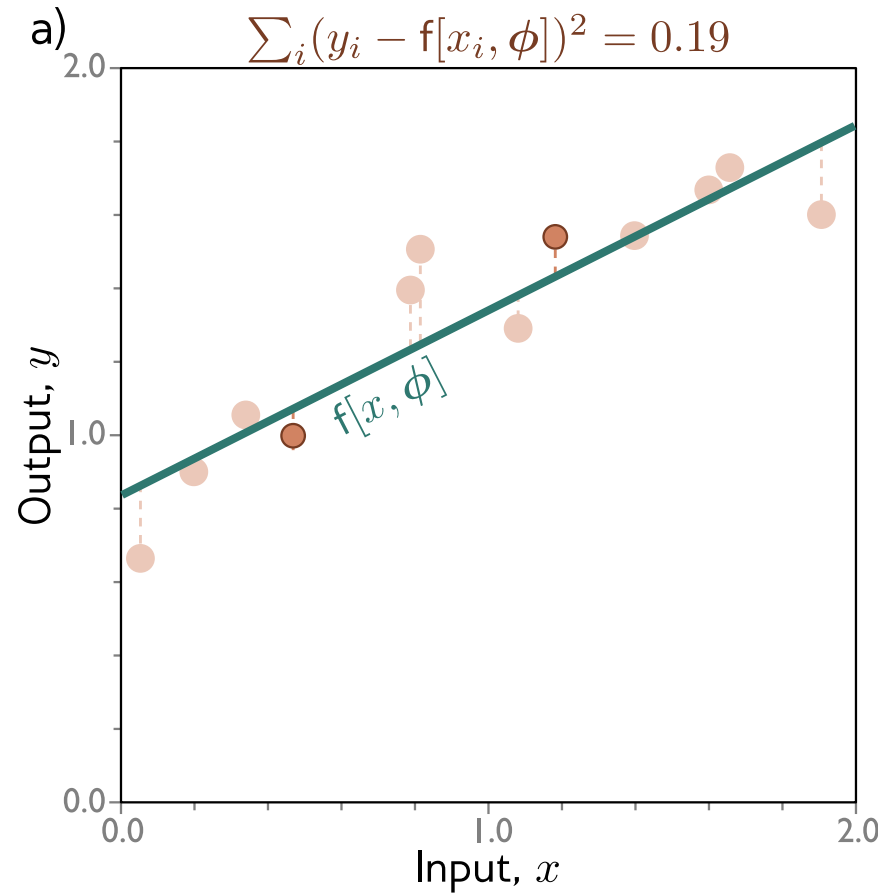$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right]$$

$$= \underset{\phi}{\operatorname{argmin}} \left[ \sum_{i=1}^{I} (y_i - \mathrm{f}[\mathbf{x}_i, \phi])^2 \right]$$

# Estimating variance

- But we could learn it during training:

$$\hat{\boldsymbol{\phi}}, \hat{\sigma}^2 = \underset{\boldsymbol{\phi}, \sigma^2}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mathrm{f}[\mathbf{x}_i, \boldsymbol{\phi}])^2}{2\sigma^2} \right] \right] \right]$$

- Do gradient descent on both model parameters, $\phi$, and the variance, $\sigma^2$

$$\frac{\partial L}{\partial \phi} \ \text{and} \ \frac{\partial L}{\partial \sigma^2}$$

# Heteroscedastic regression

- We were assuming that the noise $\sigma^2$ is the same everywhere (homoscedastic).

- But we could make the noise a function of the data *x*.

- Build a model with two outputs:

$$\mu = \mathrm{f}_1[\mathbf{x}, \phi]$$

$$\sigma^2 = \mathrm{f}_2[\mathbf{x}, \phi]^2$$

$$\hat{\phi} = \underset{\phi}{\mathrm{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ \frac{1}{\sqrt{2\pi \mathrm{f}_2[\mathbf{x}_i, \phi]^2}} \right] - \frac{(y_i - \mathrm{f}_1[\mathbf{x}_i, \phi])^2}{2\mathrm{f}_2[\mathbf{x}_i, \phi]^2} \right]$$

# Heteroscedastic regression

# Example 1: Univariate Regression Takeaways

- *Least squares loss* is a good choice assuming normal distribution
- The best prediction is the predicted mean
- We can also estimate global or local variance

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Example 2: binary classification



Real world input    Model input    Model    Model output    Real world output

"The steak was terrible, the salad was rotten, and the soup tasted like socks"

$$\begin{bmatrix} 8672 \\ 8194 \\ 9804 \\ 8634 \\ 8672 \\ \vdots \end{bmatrix}$$

Supervised learning model

$$\begin{bmatrix} 0.02 \\ 0.98 \end{bmatrix}$$

Positive
Negative

- Goal: predict which of two classes $y \in \{0, 1\}$ the input *x* belongs to

# Example 2: binary classification

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

- Domain: $y \in \{0, 1\}$

- Bernoulli distribution

- One parameter $\lambda \in [0,1]$

$$Pr(y|\lambda) = \begin{cases} 1 - \lambda & y = 0 \\ \lambda & y = 1 \end{cases}$$

$$Pr(y|\lambda) = (1 - \lambda)^{1-y} \cdot \lambda^y$$

# Example 2: binary classification

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$.

Problem:

- Output of neural network can be anything
- Parameter $\lambda \in [0,1]$

Solution:

- Pass through function that maps "anything" to [0,1]

# Example 2: binary classification

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$.

Problem:

- Output of neural network can be anything

- Parameter $\lambda \in [0,1]$

Solution:

- Pass through logistic sigmoid function that maps "anything to [0,1]:

$$\text{sig}[z] = \frac{1}{1 + \exp[-z]}$$

# Example 2: binary classification

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$.

$$Pr(y|\lambda) = (1 - \lambda)^{1-y} \cdot \lambda^y$$

$$Pr(y|\mathbf{x}) = (1 - \text{sig}[\mathbf{f}[\mathbf{x}|\boldsymbol{\phi}]])^{1-y} \cdot \text{sig}[\mathbf{f}[\mathbf{x}|\boldsymbol{\phi}]]^y$$

# Example 2: binary classification

a)



b)

$$\text{sig}\left[f[x,\phi]\right]$$

c)

# Example 2: binary classification

3. To train the model, find the network parameters $\hat{\phi}$ that minimize the negative log-likelihood loss function over the training dataset pairs $\{\mathbf{x}_i, \mathbf{y}_i\}$:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[ L[\phi] \right] = \underset{\phi}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]. \tag{5.7}$$

$$Pr(y|\mathbf{x}) = (1 - \operatorname{sig}[\mathbf{f}[\mathbf{x}|\phi]])^{1-y} \cdot \operatorname{sig}[\mathbf{f}[\mathbf{x}|\phi]]^{y}$$

$$L[\phi] = \sum_{i=1}^{I} -(1 - y_i) \log \left[ 1 - \operatorname{sig}[\mathbf{f}[\mathbf{x}_i|\phi]] \right] - y_i \log \left[ \operatorname{sig}[\mathbf{f}[\mathbf{x}_i|\phi]] \right]$$

*Binary cross-entropy loss*

# Example 2: binary classification

4. To perform inference for a new test example $\mathbf{x}$, return either the full distribution $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x},\hat{\phi}])$ or the maximum of this distribution.



Choose y=1 where $\lambda$ is greater than 0.5, otherwise 0
And we get a probability estimate!

# Example 2: Binary Classification Takeaways

- Binary cross entropy loss as the loss function

- Threshold to get prediction

- We also get a probability or "confidence value"

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1:  univariate regression
- Example 2:  binary classification
- Example 3:  multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Example 3: multiclass classification



Goal: predict which of K classes $y \in \{1, 2, \ldots, K\}$ the input $x$ belongs to

# Example 3: multiclass classification

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

- Domain: $y \in \{1, 2, \ldots, K\}$
- Categorical distribution
- K parameters $\lambda_k \in [0,1]$
- Sum of all parameters = 1

$$Pr(y = k) = \lambda_k$$

# Example 3: multiclass classification

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$.

Problem:

- Output of neural network can be anything
- Parameters $\lambda_k \in [0,1]$, sum to one

Solution:

$$\text{softmax}_k[\mathbf{z}] = \frac{\exp[z_k]}{\sum_{k'=1}^{K} \exp[z_{k'}]}$$

- Pass through function that maps "anything" to [0,1], sum to one

$$Pr(y = k|\mathbf{x}) = \text{softmax}_k[\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]]$$

# Example 3: multiclass classification



$$Pr(y = k|\mathbf{x}) = \text{softmax}_k[\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]]$$

# Example 3: multiclass classification

3. To train the model, find the network parameters $\hat{\phi}$ that minimize the negative log-likelihood loss function over the training dataset pairs $\{\mathbf{x}_i, \mathbf{y}_i\}$:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[ L[\phi] \right] = \underset{\phi}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]. \qquad (5.7)$$

$$L[\phi] = -\sum_{i=1}^{I} \log \left[ \operatorname{softmax}_{y_i} \left[ \mathbf{f}\left[ \mathbf{x}_i, \phi \right] \right] \right]$$

$$\operatorname{softmax}_k[\mathbf{z}] = \frac{\exp[z_k]}{\sum_{k'=1}^{K} \exp[z_{k'}]}$$

$$= -\sum_{i=1}^{I} \mathrm{f}_{y_i} \left[ \mathbf{x}_i, \phi \right] - \log \left[ \sum_{k=1}^{K} \exp\left[ \mathrm{f}_k \left[ \mathbf{x}_i, \phi \right] \right] \right]$$

*Multiclass cross-entropy loss*

# Example 3: multiclass classification

4. To perform inference for a new test example $\mathbf{x}$, return either the full distribution $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \hat{\boldsymbol{\phi}}])$ or the maximum of this distribution.



Choose the class with the largest probability
We also get probability or "confidence"

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1:  univariate regression
- Example 2:  binary classification
- Example 3:  multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Other data types

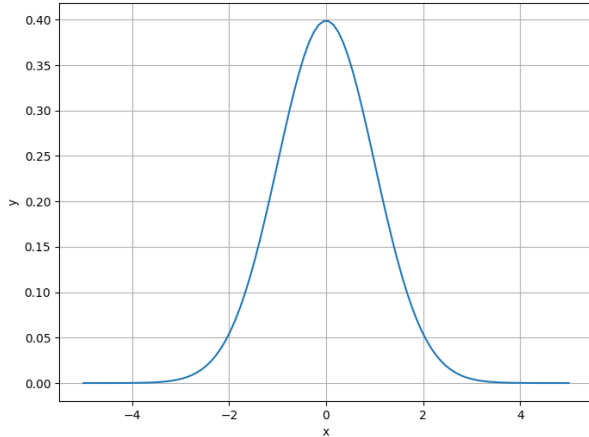| Data Type | Domain | Distribution | Use |
|---|---|---|---|
| univariate, continuous, unbounded | $y \in \mathbb{R}$ | univariate normal | regression |
| univariate, continuous, unbounded | $y \in \mathbb{R}$ | Laplace or t-distribution | robust regression |
| univariate, continuous, unbounded | $y \in \mathbb{R}$ | mixture of Gaussians | multimodal regression |
| univariate, continuous, bounded below | $y \in \mathbb{R}^+$ | exponential or gamma | predicting magnitude |
| univariate, continuous, bounded | $y \in [0, 1]$ | beta | predicting proportions |
| multivariate, continuous, unbounded | $\mathbf{y} \in \mathbb{R}^K$ | multivariate normal | multivariate regression |
| univariate, continuous, circular | $y \in (-\pi, \pi]$ | von Mises | predicting direction |
| univariate, discrete, binary | $y \in \{0, 1\}$ | Bernoulli | binary classification |
| univariate, discrete, bounded | $y \in \{1, 2, \ldots, K\}$ | categorical | multiclass classification |
| univariate, discrete, bounded below | $y \in [0, 1, 2, 3, \ldots]$ | Poisson | predicting event counts |
| multivariate, discrete, permutation | $\mathbf{y} \in \mathrm{Perm}[1, 2, \ldots, K]$ | Plackett-Luce | ranking |

**Figure 5.11** Distributions for loss functions for different prediction types.

# Other Distributions

### Gaussian



$y \in \mathbb{R}$ Regression

### Laplace



$y \in \mathbb{R}$ Robust Regression

### Mixture of Gaussians



$y \in \mathbb{R}$ Multimodal Regression

### Gamma



$y \in \mathbb{R}^+$ Predict Magnitude

### Beta



$y \in [0,1]$ Predict Proportions

### Von Mises



$y \in (-\pi, \pi]$ Predict Directions

### Poisson



$y \in [0,1,2, \dots ]$ Predict Event Counts

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Multiple outputs

- Treat each output $y_d$ as independent:

$$Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}_i, \phi]) = \prod_d Pr(y_d|\mathbf{f}_d[\mathbf{x}_i, \phi])$$

- Negative log likelihood becomes sum of terms:

$$L[\phi] = -\sum_{i=1}^{I} \log\Big[Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}_i, \phi])\Big] = -\sum_{i=1}^{I}\sum_d \log\Big[Pr(y_{id}|\mathbf{f}_d[\mathbf{x}_i, \phi])\Big]$$

# Example 4: multivariate regression

Real world input



Model input

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ 17 \\ 1 \\ 1 \\ \vdots \end{bmatrix}$$

Model

Supervised learning model

Model output

$$\begin{bmatrix} -12.9 \\ 56.4 \end{bmatrix}$$

Real world output

Freezing point is -12.9°C
Boiling point is 56.4°C

# Example 4: multivariate regression

- Goal: to predict a multivariate target $\mathbf{y} \in \mathbb{R}^{D_o}$
- Solution treat each dimension independently

$$Pr(\mathbf{y}|\boldsymbol{\mu}, \sigma^2) = \prod_{d=1}^{D_o} Pr(y_d|\mu_d, \sigma^2)$$

$$= \prod_{d=1}^{D_o} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_d - \mu_d)^2}{2\sigma^2}\right]$$

- Make network with $D_o$ outputs to predict means

$$Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}], \sigma^2) = \prod_{d=1}^{D_o} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_d - \mathrm{f}_d[\mathbf{x}, \boldsymbol{\phi}])^2}{2\sigma^2}\right]$$

# Example 4: multivariate regression

- What if the outputs vary in magnitude
  - E.g., predict weight in kilos and height in meters
  - One dimension has much bigger numbers than others
- Could learn a separate variance for each…
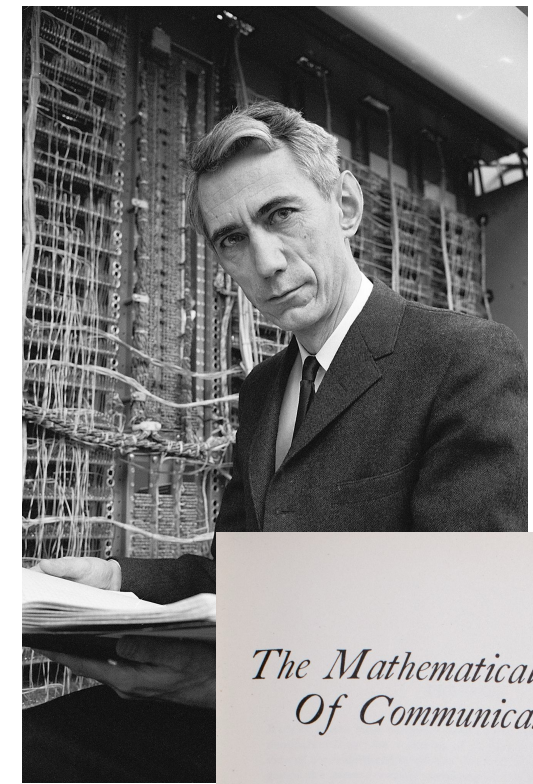- …or rescale before training, and then rescale output in opposite way

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Information Theory and Entropy



- **Claude Shannon:** the "father of information theory," was an American mathematician, electrical engineer, and cryptographer

- **Theory of Communication:** In his landmark 1948 paper, "A Mathematical Theory of Communication," Shannon introduced a formal framework for the transmission, processing, and storage of information.

- **Information Theory:** Quantified information, allowing for the measurement of information content in messages, which is crucial for data compression, error detection and correction, and more.

- **Concept of Information Entropy:** introduced entropy as a measure of the uncertainty or randomness in a set of possible messages, providing a limit on the best possible lossless compression of any communication.

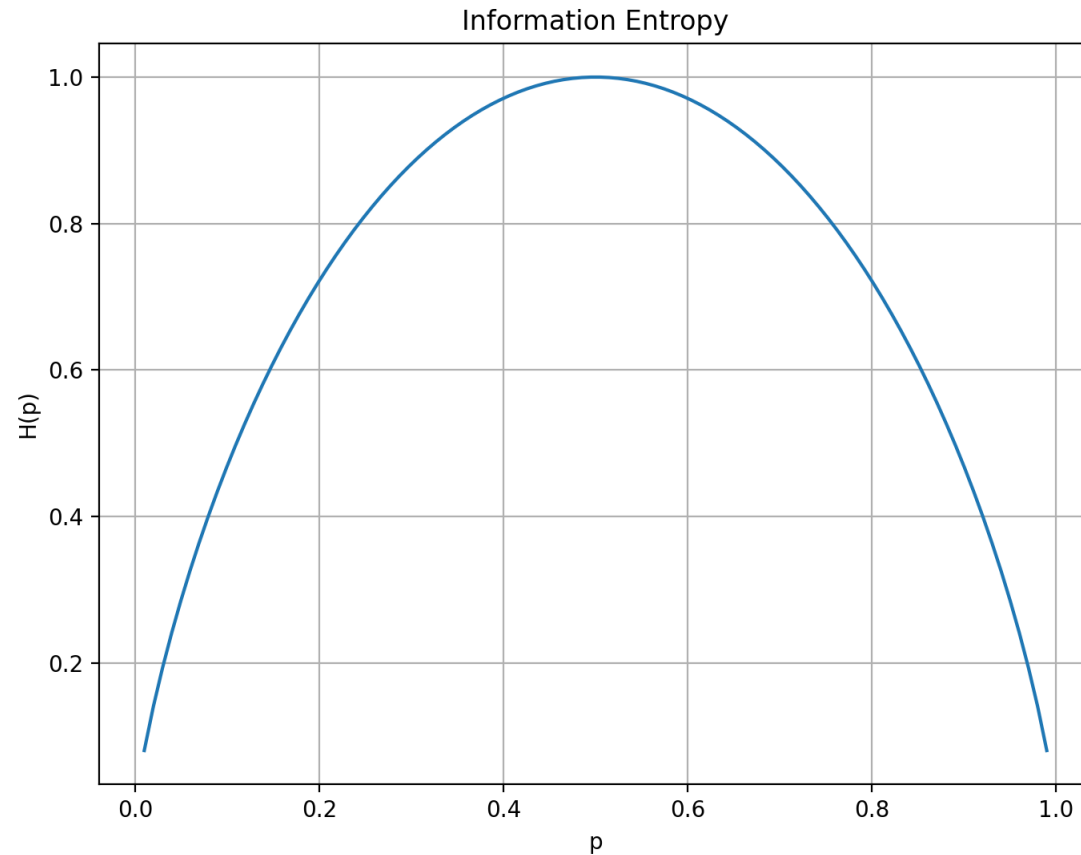$$H(x) = -\sum_{x} P(x) \log_2(P(x))$$

# Entropy for a Binary Event $\quad x \in \{0,1\}$



Information Entropy

$$H(x) = -\sum_x P(x) \log_2\big(P(x)\big) = -p \log_2(p) - (1-p)\log_2(1-p)$$

# Cross Entropy – Concept from Information Theory

Measures the difference between two probability distributions: the true distribution of the labels and the predicted distribution of the labels by a model.



a) Empirical data distribution

b) Model distribution

Minimize KL divergence

$$\text{KL}[q||p] = \int_{-\infty}^{\infty} q(z) \log[q(z)]dz - \int_{-\infty}^{\infty} q(z) \log[p(z)]dz$$

Kullback-Leibler Divergence -- a measure between probability distributions[108]

# Cross Entropy – Concept from Information Theory

For discrete distributions, the cross-entropy between two distributions $p$ and $q$ over the same underlying set of events is defined as:

$$H(p, q) = -\sum p(x) \, log \, q(x)$$

Here, $p(x)$ is the true probability of an event $x$, and $q(x)$ is the estimated probability of the same event according to the model.

For instance, in binary classification:

$$H(p, q) = -[y \log(\hat{y}) + (1 - y)\log(1 - \hat{y})$$

Here, $y$ is the true label (0 or 1), and $\hat{y}$ is the predicted probability of the class being 1.

# Recap

- Reconsidered loss functions as fitting a parametric probability model
- Introduced Maximum Likelihood criterion for finding parameters to making the training data most probably under that model
- Introduced a 4-step recipe for (1) picking a suitable parametric probability distribution, (2) defining the model to pick one or more of the parameters, (3) training the model and (4) doing inference
- Derived loss functions for univariate regression, binary and multiclass classification
- Briefly reviewed parametric probability models for other types of data
- Discussed how this is the same as Cross Entropy from Information Theory

# Next up

- Now let's find the parameters that give the smallest loss
  - ➔ Training the model

# Feedback?