



ESTIMATION AND REGRESSION

EE 541 – UNIT 3B



REGRESSION OVERVIEW

- Regression is data fitting to a specific parameterized function class
- Linear regression
 - Same as LMMSE, but with data averages replacing expectation (ensemble averages)
 - Linear least-squares
 - Generalize on-line learning to full-batch and mini-batches
- Regularization (*later*)
- Logistical Regression (*later*)



GENERAL REGRESSION PROBLEM

Given a data set: $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$

General regression problem:

$$\min_{\Theta} \langle C(\mathbf{y}, \mathbf{g}(\mathbf{x}; \Theta)) \rangle_{\mathcal{D}} \quad \Theta_{opt} = \arg \min_{\Theta} \langle C(\mathbf{y}, \mathbf{g}(\mathbf{x}; \Theta)) \rangle_{\mathcal{D}} \quad \hat{\mathbf{y}} = \mathbf{g}(\mathbf{x}; \Theta_{opt})$$

Empirical expectation (average over data):

$$\langle \mathbf{h}(\mathbf{x}, \mathbf{y}) \rangle_{\mathcal{S}} \equiv \frac{1}{|\mathcal{S}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathcal{S}} \mathbf{h}(\mathbf{x}_n, \mathbf{y}_n)$$

x ~ regressor (observed)
y ~ target (desired)

For large averaging sets (i.e., many realizations):

$$\mathbb{E}\{\mathbf{h}(\mathbf{x}(t), \mathbf{y}(t))\} = \int \mathbf{h}(\mathbf{x}, \mathbf{y}) p_{\mathbf{x}(t), \mathbf{y}(t)}(\mathbf{x}, \mathbf{y}) d\mathbf{y} d\mathbf{x} \approx \langle \mathbf{h}(\mathbf{x}, \mathbf{y}) \rangle_{\mathcal{S}} \quad \text{sample mean}$$

Monte Carlo method



LEAST-SQUARES (LS) REGRESSION PROBLEM

$$\min_{\Theta} \langle \|\mathbf{y} - \mathbf{g}(\mathbf{x}; \Theta)\| \rangle_{\mathcal{D}} \quad \Leftrightarrow \quad \min_{\Theta} \sum_{n=1}^N \|\mathbf{y}_n - \mathbf{g}(\mathbf{x}_n; \Theta)\|^2$$

$$\Theta_{\text{opt}} = \arg \min_{\Theta} \langle \|\mathbf{y} - \mathbf{g}(\mathbf{x}; \Theta)\|^2 \rangle_{\mathcal{D}}$$

Squared-error is a common cost function in (electrical) engineering

corresponds to **power or energy** in many applications



LINEAR AND AFFINE LEAST SQUARES REGRESSION

Linear regression problem:

$$\min_{\mathbf{W}} \langle \|\mathbf{y} - \mathbf{W}\mathbf{x}\|^2 \rangle_{\mathcal{D}} \Leftrightarrow \min_{\mathbf{W}} \sum_{n=1}^N \|\mathbf{y}_n - \mathbf{W}\mathbf{x}_n\|^2$$

$$\mathbf{W}_{\text{LLSE}} = \arg \min_{\mathbf{W}} \langle \|\mathbf{y} - \mathbf{W}\mathbf{x}\|^2 \rangle_{\mathcal{D}}$$

$$\hat{\mathbf{y}} = \mathbf{W}_{\text{LLSE}}\mathbf{x}$$

Affine regression (*a.k.a.*, Linear regression):

$$\mathbf{W}_{\text{ALSE}}, \mathbf{b}_{\text{ALSE}} = \arg \min_{\mathbf{W}, \mathbf{b}} \langle \|\mathbf{y} - [\mathbf{W}\mathbf{x} + \mathbf{b}]\|^2 \rangle_{\mathcal{D}}$$

$$\hat{\mathbf{y}} = \mathbf{W}_{\text{ALSE}}\mathbf{x} + \mathbf{b}_{\text{ALSE}}$$



LINEAR AND AFFINE REGRESSION SOLUTION

Data averaging operator has linearity property like expectation

$$\mathbb{E}\{L(x(t))\} = L(\mathbb{E}(x(t))) \qquad \langle L(x) \rangle = L(\langle x \rangle)$$

This means the solutions are the same as the MMSE solutions with expectation replaced by data average

For example, Linear LS regression:

$$\mathbf{W}_{LLSE} = \hat{\mathbf{R}}_{YX} \hat{\mathbf{R}}_X^{-1}$$

$$\hat{\mathbf{y}} = \hat{\mathbf{R}}_{YX} \hat{\mathbf{R}}_X^{-1} \mathbf{x}$$

$$\begin{aligned} LLS\mathcal{E} &= \langle \|\mathbf{y} - \mathbf{W}_{LLSE} \mathbf{x}\|^2 \rangle \\ &= \langle \|\mathbf{y}\|^2 - \|\mathbf{W}_{LLSE} \mathbf{x}\|^2 \rangle_{\mathcal{D}} \\ &= \text{Tr}(\hat{\mathbf{R}}_Y - \hat{\mathbf{R}}_{YX} \hat{\mathbf{R}}_X^{-1} \hat{\mathbf{R}}_{XY}) \end{aligned}$$

$$\begin{aligned} \hat{\mathbf{R}}_X &= \langle \mathbf{x} \mathbf{x}^T \rangle_{\mathcal{D}} \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \end{aligned}$$

$$\begin{aligned} \hat{\mathbf{R}}_{XY} &= \langle \mathbf{x} \mathbf{y}^T \rangle_{\mathcal{D}} \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{y}_n^T \end{aligned}$$

PROOF FOR LLSE REGRESSION

$$\min_{\mathbf{W}} \langle \|\mathbf{y} - \mathbf{W}\mathbf{x}\|^2 \rangle_{\mathcal{D}}$$

Proof:

$$\text{LSE}(\mathbf{G}) = \langle \|\mathbf{y} - \mathbf{G}\mathbf{x}\|^2 \rangle$$

$$= \langle \|\mathbf{y} - \mathbf{G}_{\text{opt}}\mathbf{x} + (\mathbf{G}_{\text{opt}} - \mathbf{G})\mathbf{x}\|^2 \rangle$$

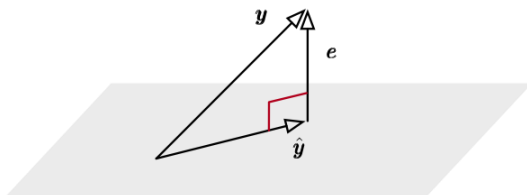
$$= \langle \|\mathbf{y} - \mathbf{G}_{\text{opt}}\mathbf{x}\|^2 \rangle + \text{Tr} \left((\mathbf{G}_{\text{opt}} - \mathbf{G}) \hat{\mathbf{R}}_{\mathbf{x}} (\mathbf{G}_{\text{opt}} - \mathbf{G})^T \right)$$

$$+ 2 \text{Tr} \left((\hat{\mathbf{R}}_{\mathbf{y}\mathbf{x}} - \mathbf{G}_{\text{opt}} \hat{\mathbf{R}}_{\mathbf{x}}) (\mathbf{G}_{\text{opt}} - \mathbf{G})^T \right)$$

$$\mathbf{v}^T \mathbf{w} = \text{Tr}(\mathbf{w} \mathbf{v}^T)$$

if: $\mathbf{G}_{\text{opt}} \hat{\mathbf{R}}_{\mathbf{x}} = \hat{\mathbf{R}}_{\mathbf{y}\mathbf{x}}$ then: $\text{LSE}(\mathbf{G}) = \mathbb{E} \left\{ \|\mathbf{y} - \mathbf{G}_{\text{opt}}\mathbf{x}\|^2 \right\} + \underbrace{\text{Tr} \left((\mathbf{G}_{\text{opt}} - \mathbf{G}) \hat{\mathbf{R}}_{\mathbf{x}} (\mathbf{G}_{\text{opt}} - \mathbf{G})^T \right)}_{\geq 0 \forall \mathbf{G}, \text{ since } \hat{\mathbf{R}}_{\mathbf{x}} \text{ is psd}}$

Wiener-Hopf equations (*Orthogonality Principle*)



space of all estimates/approximations

because of orthogonality principle
(error and signal uncorrelated)

$$\begin{aligned} \langle \|\mathbf{y} - \hat{\mathbf{y}}\|^2 \rangle &= \langle \|\mathbf{y}\|^2 \rangle + \langle \|\hat{\mathbf{y}}\|^2 \rangle \\ &= \text{Tr}(\hat{\mathbf{R}}_{\mathbf{y}} - \hat{\mathbf{R}}_{\mathbf{y}\mathbf{x}} \hat{\mathbf{R}}_{\mathbf{x}}^{-1} \hat{\mathbf{R}}_{\mathbf{x}\mathbf{y}}) \end{aligned}$$



SOLUTION TO LINEAR AND AFFINE (LS) REGRESSION

It makes intuitive sense:

For LMMSE estimate: if you did not know the second moments you would estimate these correlations from data

in addition to optimality in the Gaussian case, linear MMSE estimation is popular because it requires much less data to accurately estimate second moments than a complete statistical description (or higher moments)

LLSE REGRESSION: SCALAR FROM SCALAR

Estimate y from x

Linear regression problem:

$$\min_w \langle (y - wx)^2 \rangle \Leftrightarrow \min_w \frac{1}{N} \sum_{n=1}^N (y_n - wx_n)^2$$

Solution (special case):

$$w_{LLSE} = \frac{\hat{r}_{yx}}{\hat{r}_x}$$

$$\hat{y} = \frac{\hat{r}_{yx}}{\hat{r}_x} x$$

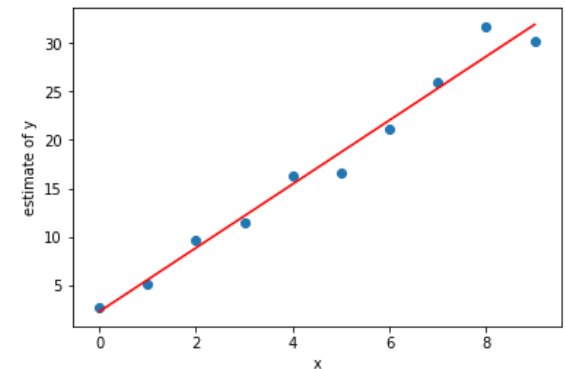
$$\hat{r}_x = \langle x^2 \rangle = \frac{1}{N} \sum_{n=1}^N x_n^2$$

$$\hat{r}_{yx} = \langle yx \rangle = \frac{1}{N} \sum_{n=1}^N y_n x_n$$

$$\begin{aligned} LLS\epsilon &= \langle [y - w_{LLSE}x]^2 \rangle \\ &= \langle y^2 \rangle - \langle [w_{LLSE}x]^2 \rangle \\ &= \hat{r}_y - \hat{r}_{yx}^2 \hat{r}_x^{-1} \end{aligned}$$

if sample means all 0:

$$= \hat{\sigma}_y^2 (1 - \hat{\rho}^2)$$



LLSE REGRESSION: SCALAR FROM SCALAR

Estimate y from x

Linear regression problem:

$$\min_w \langle (y - wx)^2 \rangle \Leftrightarrow \min_w \frac{1}{N} \sum_{n=1}^N (y_n - wx_n)^2 \Leftrightarrow \|y - wx\|^2$$

Solution (special case):

$$w_{LLSE} = \frac{\mathbf{y}^T \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

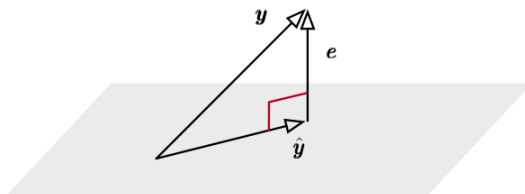
$$(N)LLSE = \|\mathbf{y}\|^2 - \left(\frac{\mathbf{y}^T \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \right)^2 \|\mathbf{x}\|^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

$$\hat{\mathbf{y}} = \frac{\mathbf{y}^T \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \mathbf{x}$$

$$= \|\mathbf{y}\|^2 - \frac{(\mathbf{y}^T \mathbf{x})^2}{\|\mathbf{x}\|^2}$$



space of all estimates/approximations

this “stacked” approach yields the same as the $\langle \cdot \rangle_D$ approach on the previous slides

$\hat{\mathbf{y}}$ stacked in a vector

LLSE REGRESSION: SCALAR FROM VECTOR

Estimate y from \mathbf{x}

Linear regression problem:

$$\min_{\mathbf{w}} \left\langle (y - \mathbf{w}^T \mathbf{x})^2 \right\rangle \Leftrightarrow \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2$$

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{w} = \hat{\mathbf{R}}_X^{-1} \hat{\mathbf{r}}_{xy}$$

$$\hat{\mathbf{r}}_{xy} = \hat{\mathbf{R}}_{xy} = \langle \mathbf{x} y \rangle$$

$$\hat{\mathbf{R}}_X \mathbf{w} = \hat{\mathbf{r}}_{xy} \quad \text{“Normal Equations”}$$

$$LLS\varepsilon = \hat{r}_y - \hat{\mathbf{r}}_{xy}^T \hat{\mathbf{R}}_X^{-1} \hat{\mathbf{r}}_{xy}$$

$$\hat{\mathbf{R}}_X \mathbf{w} = \hat{\mathbf{r}}_{xy}$$

similar: just change $\mathbb{E}[\cdot]$ to $\langle \cdot \rangle_D$ in LMMSE result

what about the “stacked” approach for this case??

LLSE REGRESSION: SCALAR FROM VECTOR

Estimate y from x

Linear regression problem:

$$\min_w \left\langle (y - w^T x)^2 \right\rangle \Leftrightarrow \min_w \frac{1}{N} \sum_{n=1}^N (y_n - w^T x)^2 \Leftrightarrow \min_w \|y - Xw\|^2$$

Solution (special case):

$$\hat{y} = Xw$$

$$w = (X^T X)^{-1} X^T y$$

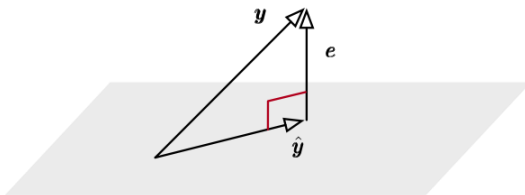
$$\begin{aligned} \hat{y} &= X(X^T X)^{-1} X^T y \\ &= P_X y \end{aligned}$$

$$\begin{aligned} (N)LLS\varepsilon &= \text{Tr}(\|y\|^2 - \|P_X y\|^2) \\ &= \text{Tr}(\|(I - P_X)y\|^2) \end{aligned}$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} \quad X^T = [x_1 \quad x_2 \quad \cdots \quad x_N]$$

$$X^T X = X^T y$$

normal equations



space of all estimates/approximations

this is the same as $\langle \cdot \rangle_D$ case, with all \hat{y} stacked in a vector

$$\hat{R}_X^{-1} \hat{r}_{xy} = \left(\frac{1}{N} X^T X \right)^{-1} \left[\frac{1}{N} X^T y \right]$$

THE AFFINE TO LINEAR MATH “TRICK”

$$\min_{\mathbf{w}} \left\langle (y - \mathbf{w}^T \mathbf{x})^2 \right\rangle \Leftrightarrow \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2$$

$$\begin{aligned} \hat{y} &= [\mathbf{x}^T \quad | \quad 1] \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \\ &= \mathbf{x}^T \mathbf{w} + b \end{aligned}$$

$$\begin{aligned} \hat{y} &= [\mathbf{X} \quad | \quad \mathbf{1}] \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \\ &= \mathbf{X}\mathbf{w} + \mathbf{b}1 \end{aligned}$$

$$[\mathbf{X} \quad | \quad \mathbf{1}] = \left[\begin{array}{c|c} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_N^T & 1 \end{array} \right]$$

therefore: compact notation even if using bias (b) term



LINEAR CLASSIFICATION



LINEAR CLASSIFIER

perform linear regression and then **threshold** to hard decision

Example:

$$y \in \{-1, +1\}$$

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

$$\text{sign}(v) = \begin{cases} +1, & v \geq 0 \\ -1, & v < 0 \end{cases}$$

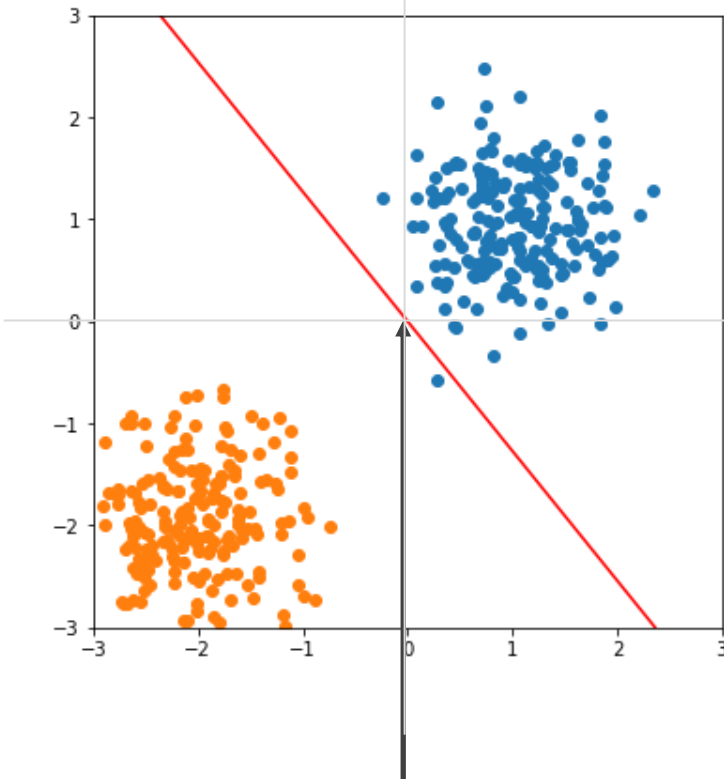
$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2$$

standard LLSE regression with prediction thresholding

Review: [linear_classifier_examples.ipynb](#)

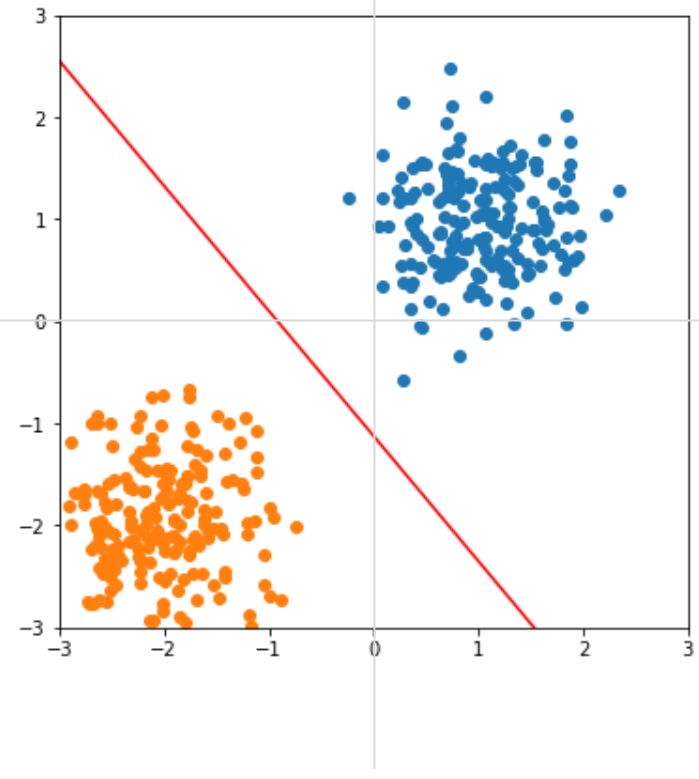
EXAMPLE: LINEAR AND AFFINE REGRESSION

$$\hat{y} = \mathbf{x}^T \mathbf{w}$$



for the case with no bias term, the decision threshold must pass through the origin

$$\hat{y} = [\mathbf{x}^T \mid 1] \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$



adding the bias term allows for offset from the origin

$$s_0 = \begin{bmatrix} +1 \\ +1 \end{bmatrix}$$

$$s_1 = \begin{bmatrix} -2 \\ -2 \end{bmatrix}$$



MAXIMUM LIKELIHOOD ESTIMATION EXAMPLE

this is a model for the data $\{(x_n, y_n)\}$:

$$y_n = \mathbf{w}^T \mathbf{x}_n + v_n, \quad n = 1, 2, \dots, N$$

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \mathbf{v}(t)$$

$$p_{v(t)}(\mathbf{v}) = \mathcal{N}_N(\mathbf{v}; \mathbf{0}, \sigma_v^2 \mathbf{I})$$

$$p_{\mathbf{y}(t)|\mathbf{X}(t)}(\mathbf{y}|\mathbf{X}; \mathbf{w}) = p_{\mathbf{v}(t)}(\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathcal{N}_N(\mathbf{v}; \mathbf{X}\mathbf{w}, \sigma_v^2 \mathbf{I})$$

$$NLL(\mathbf{w}) = -\ln(p_{\mathbf{y}(t)}(\mathbf{y}|\mathbf{X}; \mathbf{w}))$$

$$= -\ln\left(\frac{1}{(2\pi\sigma_v^2)^{\frac{N}{2}}} \exp\left[-\frac{1}{2\sigma_v^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2\right]\right)$$

$$= -\frac{1}{2\sigma_v^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \frac{N}{2} \ln(2\pi\sigma_v^2)$$

$$\max_{\mathbf{w}} p_{\mathbf{y}(t)|\mathbf{X}(t)}(\mathbf{y}|\mathbf{X}; \mathbf{w}) \Leftrightarrow \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

Maximum Likelihood \Leftrightarrow Minimize Neg-Log-Likelihood \Leftrightarrow LLSE regression

(under this model for the data)



PROPERTIES OF ML ESTIMATORS

- Asymptotically Gaussian:
 - For large amounts of data, the ML estimate is Gaussian with mean equal to the true parameter (models matched)
- Consistent:
 - The limit in probability of the ML estimate is the true parameter (model matched)
- The ML estimate minimizes the KL Divergence between the model distribution and the empirical data distribution. KL divergence measures the difference between two distribution (Info. Theory).
 - Minimizing KL divergence in this case also corresponds to minimizing the cross entropy



INFORMATION THEORY



ML ESTIMATION AND INFORMATION THEORY

Entropy:

$$\begin{aligned} H(X(t)) &= \mathbb{E} \left\{ \log_2 \left(\frac{1}{p_{X(t)}(X(t))} \right) \right\} \\ &= \sum_k p_{X(t)}(k) \log_2 \left(\frac{1}{p_{X(t)}(k)} \right) \\ &= \sum_k p_k \log_2 \left(\frac{1}{p_k} \right) \end{aligned}$$

Intuition:

events with low probability have large information —
e.g., “it will snow in Phoenix tomorrow”

the entropy is the average information learned when the
value of $X(u)$ is revealed.

Examples:

weather report in Phoenix has low entropy (almost always the same), whereas
in Sioux City, SD it has high entropy (highly variant weather)

fair die:

$$H(X(t)) = \log_2(1/6) = 2.58 \text{ bits/roll}$$

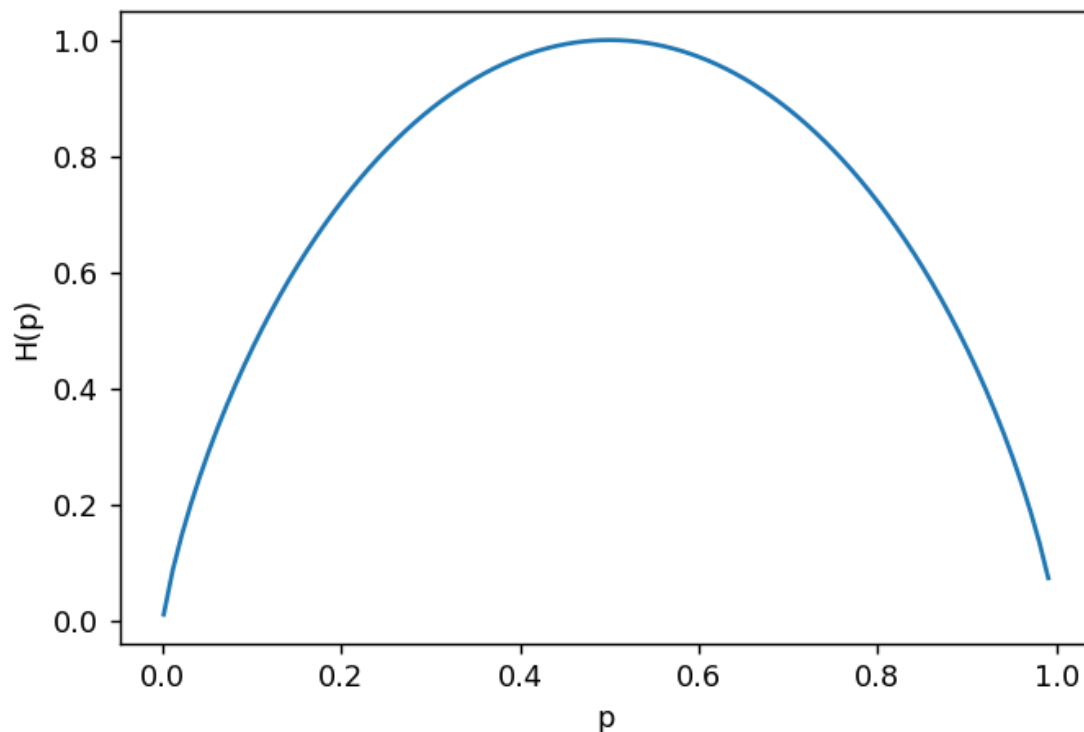
loaded die:

$$\begin{aligned} H(X(t)) &= -0.4 \log_2(0.4) - 0.1 \log_2(0.1) - 0.01 \log_2(0.01) \\ &\quad - 0.09 \log_2(0.09) - 0.25 \log_2(0.25) - 0.15 \log_2(0.15) \\ &= 2.15 \text{ bits/roll} \end{aligned}$$

ML ESTIMATION AND INFORMATION THEORY

Entropy of i.i.d. Bernoulli Source (with success probability p)

$$H(p) = -p \log_2(p) - (1 - p) \log_2(1 - p)$$





ML ESTIMATION AND INFORMATION THEORY

KL-Divergence

$$\begin{aligned} D(p \parallel \tilde{p}) &= \mathbb{E}_p \left\{ \log \left(\frac{p_x(X(t))}{\tilde{p}_x(X(t))} \right) \right\} \\ &= \sum_k p_k \log \left(\frac{p_k}{\tilde{p}_k} \right) \\ &= \sum_k p_k \log(p_k) - \sum_k p_k \log(\tilde{p}_k) \\ &= CE(p, \tilde{p}) - H(p) \end{aligned}$$

Cross-Entropy

$$CE(p, \tilde{p}) = \mathbb{E}_p \left\{ \log \left(\frac{1}{\tilde{p}(X(t))} \right) \right\}$$

ML ESTIMATION AND INFORMATION THEORY

ML parameter estimation minimizes empirical CE (and KL divergence)

$p_{data}(y|\mathbf{x})$ = data distribution of the data (typically unknown)

$p_{model}(y|\mathbf{x}; \Theta)$ = modeled distribution of the data (function of parameters)

$$\begin{aligned} CE(p_{data}, p_{model}(\Theta)) &= \mathbb{E}_{p_{data}(y|x)} \left\{ \log \left(\frac{1}{p_{model}(y(t)|\mathbf{x}(t); \Theta)} \right) \right\} \\ &\approx \langle -\log(p_{model}(y|\mathbf{x}; \Theta)) \rangle_{\mathcal{D}} \\ &= -\frac{1}{N} \sum_{n=1}^N \log(p_{model}(y_n|\mathbf{x}_n; \Theta)) \end{aligned}$$

$$\max_{\Theta} p_{model}(\mathbf{y}|\mathbf{X}; \Theta) \Leftrightarrow \min_{\Theta} (-\log(p_{model}(\mathbf{y}|\mathbf{X}; \Theta)))$$

$$\Leftrightarrow \min_{\Theta} \left(-\sum_{n=1}^N \log(p_{model}(y_n|\mathbf{x}_n; \Theta)) \right)$$

$$\Leftrightarrow \min_{\Theta} \left(-\frac{1}{N} \sum_{n=1}^N \log(p_{model}(y_n|\mathbf{x}_n; \Theta)) \right)$$

Max-Likelihood Estimation of neural network weights is always minimizing the empirical cross entropy between data distribution and the modeled distribution

(assume i.i.d. y_n)

(empirical Cross-Entropy)



MULTI-CLASS CROSS ENTROPY EXAMPLE (“ONE HOT”)

One hot encoding: cat: 0
 dog: 1
 bird: 2

Sample data labels: $n = 1: y = 1$ (dog)
 $n = 2: y = 2$ (bird)
 $n = 3: y = 0$ (cat)

Classifier Output: $n = 1: [0.3 \quad 0.5 \quad 0.2]$ $[p(cat) \quad p(dog) \quad p(bird)]$
 $n = 2: [0.0 \quad 0.0 \quad 1.0]$
 $n = 3: [0.4 \quad 0.5 \quad 0.1]$

$$Loss = -\frac{1}{3} [\log(0.5) + \log(1.0) + \log(0.4)]$$

$$\overline{MCE}(p_{data}, p_{model}(w)) = -\frac{1}{N} \sum_{n=1}^N \sum_{m=1}^M \mathbb{I}[y_n = m] \log(p_{model}(y_n = m; w))$$