

COMP90051

Workshop Week 04

About the Workshops

- 7 sessions in total
 - Tue 12:00-13:00 AH211
 - Tue 12:00-13:00 AH108 *
 - Tue 13:00-14:00 AH210
 - Tue 16:15-17:15 AH109
 - Tue 17:15-18:15 AH236 *
 - Tue 18:15-19:15 AH236 *
 - Fri 14:15-15:15 AH211

About the Workshops

- Homepage

- <https://trevorcohn.github.io/comp90051-2017/workshops>

- Solutions will be released on next Friday (a week later).

Syllabus

1	Introduction; Probability theory	Probabilistic models; Parameter fitting	
2	Linear regression; Intro to regularization	Logistic regression; Basis expansion	
3	Optimization; Regularization	Perceptron	←
4	Backpropagation	CNNs; Auto-encoders	
5	Hard-margin SVMs	Soft-margin SVMs	
6	Additional topics	Kernel methods	
7	Unsupervised learning	Unsupervised learning	
8	Dimensionality reduction; Principal component analysis	Multidimensional scaling; Spectral clustering	
9	Bayesian fundamentals	Bayesian inference with conjugate priors	
10	PGMs, fundamentals	Conditional independence	
11	PGMs, inference	Belief propagation	
12	Statistical inference; Apps	Subject review	

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Supervised learning as an optimization problem
 - ❑ Perceptron update rule & loss function
 - ❑ Logistic regression
 - ❑ Predict function
 - ❑ Log loss (a.k.a. cross entropy)
- ❑ Notebook tasks
 - ❑ Task 1: Logistic regression
 - ❑ Task 2: Perceptron classifier

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Supervised learning as an optimization problem
 - ❑ Perceptron update rule & loss function
 - ❑ Logistic regression
 - ❑ Predict function
 - ❑ Log loss (a.k.a. cross entropy)
- ❑ Notebook tasks
 - ❑ Task 1: Logistic regression
 - ❑ Task 2: Perceptron classifier

Supervised learning as an optimization problem

□ Dataset

- Preprocessing / normalization / feature selection

- Split into **train***/**test**, where **test** served as the **held-out set**

 - Split **train*** into **train/validation** (or k folds **train/validation**, CV)

□ Model / Objective function

- **Parameters**, solved either analytically or by an optimizer

 - Solved on the **training set** (or **training sets** in k folds, CV)

- **Hyper-parameters**, e.g. regularization parameter

 - Selected on the **validation set** (or **validation sets** in k folds, CV)

□ **Evaluation** on the **held-out set**

Supervised learning as an optimization problem

- Dataset

 - Preprocessing / normalization / feature selection

 - Split into **train***/**test**, where **test** served as the **held-out set**

- Model / Objective function

 - **Parameters**, solved either analytically or by an optimizer

 - Solved on the **training set***

- **Evaluation** on the **held-out set**

How to solve an optimization problem?

- Analytic solution → solve it by a formula

- $\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2, \lambda \geq 0 \rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

- Iterative methods → solve it by an optimization algorithm

- To minimize an objective

- Coordinate descent

- Gradient-based optimization algorithms (optimizers)

- Simplest: gradient descent & stochastic gradient descent

- BFGS (in `4a_logistic_regression.ipynb`)

- Many more in packages...

- To maximize an objective

- Convert it to an equivalent minimization problem

-
- ❑ Linear regression has an analytic solution (with L2)
 - ❑ Perceptron has its own update rule
 - ❑ can be interpreted as using stochastic gradient descent (with an appropriate loss function defined)
 - ❑ Gradient-based optimizers can be used for
 - ❑ Linear regression
 - ❑ Support vector machines
 - ❑ Logistic regression, neural networks
 - ❑ Deep neural networks are usually constructed and optimized using special packages...

Gradient-based optimizers in

TensorFlow™

- `tf.train.Optimizer`
- `tf.train.GradientDescentOptimizer`
- `tf.train.AdadeltaOptimizer`
- `tf.train.AdagradOptimizer`
- `tf.train.AdagradDAOptimizer`
- `tf.train.MomentumOptimizer`
- `tf.train.AdamOptimizer`
- `tf.train.FtrlOptimizer`
- `tf.train.ProximalGradientDescentOptimizer`
- `tf.train.ProximalAdagradOptimizer`
- `tf.train.RMSPropOptimizer`

https://www.tensorflow.org/api_guides/python/train

Gradient-based optimizers in PYTORCH

```
class torch.optim.Optimizer(params, defaults)[source]
```

```
class torch.optim.Adadelta(params, lr=1.0, rho=0.9, eps=1e-06, weight_decay=0)[source]
```

```
class torch.optim.Adagrad(params, lr=0.01, lr_decay=0, weight_decay=0)[source]
```

```
class torch.optim.Adam(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08, weight_decay=0)[source]
```

```
class torch.optim.Adamax(params, lr=0.002, betas=(0.9, 0.999), eps=1e-08, weight_decay=0)[source]
```

```
class torch.optim.ASGD(params, lr=0.01, lambd=0.0001, alpha=0.75, t0=1000000.0, weight_decay=0)[source]
```

```
class torch.optim.LBFGS(params, lr=1, max_iter=20, max_eval=None, tolerance_grad=1e-05, tolerance_change=1e-09, history_size=100, line_search_fn=None)[source]
```

```
class torch.optim.RMSprop(params, lr=0.01, alpha=0.99, eps=1e-08, weight_decay=0, momentum=0, centered=False)[source]
```

```
class torch.optim.Rprop(params, lr=0.01, etas=(0.5, 1.2), step_sizes=(1e-06, 50))[source]
```

```
class torch.optim.SGD(params, lr=<object object>, momentum=0, dampening=0, weight_decay=0, nesterov=False)[source]
```

<http://pytorch.org/docs/master/optim.html#algorithms>

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Supervised learning as an optimization problem
 - ❑ Perceptron update rule & loss function
 - ❑ Logistic regression
 - ❑ Predict function
 - ❑ Log loss (a.k.a. cross entropy)
- ❑ Notebook tasks
 - ❑ Task 1: Logistic regression
 - ❑ Task 2: Perceptron classifier

Perceptron update rule (2-D points)

□ Data points

$$\square \mathbf{x}_1 = [1 \quad x_{1,1} \quad x_{1,2}] \rightarrow y_1 = +1, \mathbf{x}_2 = [1 \quad x_{2,1} \quad x_{2,2}] \rightarrow y_2 = +1$$

$$\square \mathbf{x}_3 = [1 \quad x_{3,1} \quad x_{3,2}] \rightarrow y_3 = -1, \mathbf{x}_4 = [1 \quad x_{4,1} \quad x_{4,2}] \rightarrow y_4 = -1$$

□ Model parameters

$$\square \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

□ Decision function

$$\square f(\mathbf{x}; \mathbf{w}) = \mathbf{x}\mathbf{w} = x_0w_0 + x_1w_1 + x_2w_2$$

$$\square \text{Predict } +1 \text{ if } f(\mathbf{x}; \mathbf{w}) > 0, \text{ predict } -1 \text{ if } f(\mathbf{x}; \mathbf{w}) < 0$$

Perceptron update rule (2-D points)

- Iterate over all the points

- Current point is \mathbf{x}_i , and its label is y_i

- Call decision function $s = f(\mathbf{x}_i; \mathbf{w}) = \mathbf{x}_i \mathbf{w}$

- Predict \hat{y}_i as $+1$ if $s > 0$, as -1 if $s < 0$

- If $y_i \neq \hat{y}_i$

- If $y_i = +1$, $\mathbf{w}^{new} = \mathbf{w}^{old} + \mathbf{x}_i^T$, or
$$\begin{cases} w_0^{new} = w_0^{old} + x_{i,0} \\ w_1^{new} = w_1^{old} + x_{i,1} \\ w_2^{new} = w_2^{old} + x_{i,2} \end{cases}$$

- If $y_i = -1$, $\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{x}_i^T$, or
$$\begin{cases} w_0^{new} = w_0^{old} - x_{i,0} \\ w_1^{new} = w_1^{old} - x_{i,1} \\ w_2^{new} = w_2^{old} - x_{i,2} \end{cases}$$

Suppose \mathbf{x}_1 is misclassified

□ This means $y_1 = +1$, $\hat{y}_1 = -1$, $s = f(\mathbf{x}_1; \mathbf{w}^{old}) = \mathbf{x}_1 \mathbf{w}^{old} < 0$

□ Apply the update rule:

□ $\mathbf{w}^{new} = \mathbf{w}^{old} + \mathbf{x}_1^T$

□ How does s change?

□ $s^{new} = \mathbf{x}_1 \mathbf{w}^{new} = \mathbf{x}_1 (\mathbf{w}^{old} + \mathbf{x}_1^T) = \mathbf{x}_1 \mathbf{w}^{old} + \mathbf{x}_1 \mathbf{x}_1^T > s^{old}$

□ We hope $s > 0$ (because $y_1 = +1$)

□ After updating \mathbf{w} , $s^{new} > s^{old}$, s^{new} may still < 0

□ But it improves a bit (at least for \mathbf{x}_1)

Suppose \mathbf{x}_3 is misclassified

□ This means $y_3 = -1$, $\hat{y}_3 = +1$, $s = f(\mathbf{x}_3; \mathbf{w}^{old}) = \mathbf{x}_3 \mathbf{w}^{old} > 0$

□ Apply the update rule:

□ $\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{x}_3^T$

□ How does s change?

□ $s^{new} = \mathbf{x}_3 \mathbf{w}^{new} = \mathbf{x}_3 (\mathbf{w}^{old} - \mathbf{x}_3^T) = \mathbf{x}_3 \mathbf{w}^{old} - \mathbf{x}_3 \mathbf{x}_3^T < s^{old}$

□ We hope $s < 0$ (because $y_3 = -1$)

□ After updating \mathbf{w} , $s^{new} < s^{old}$, s^{new} may still > 0

□ But it improves a bit (at least for \mathbf{x}_3)

A uniform update rule

□ For misclassified positive instance \mathbf{x} , $y = +1$, hope $s > 0$

□ $\mathbf{w}^{new} = \mathbf{w}^{old} + \mathbf{x}^T$ to increase $s = f(\mathbf{x}; \mathbf{w}) = \mathbf{xw}$

□ For misclassified negative instance \mathbf{x} , $y = -1$, hope $s < 0$

□ $\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{x}^T$ to decrease $s = f(\mathbf{x}; \mathbf{w}) = \mathbf{xw}$

□ A uniform update rule:

□ For misclassified instance \mathbf{x} , $\mathbf{w}^{new} = \mathbf{w}^{old} + y\mathbf{x}^T$ so that

□ $s = f(\mathbf{x}; \mathbf{w}) = \mathbf{xw}$ is increased if $y = +1$

□ $s = f(\mathbf{x}; \mathbf{w}) = \mathbf{xw}$ is decreased if $y = -1$

Perceptron update rule (2-D points)

- Iterate over all the points
 - Current point is \mathbf{x}_i , and its label is y_i
 - Call decision function $s = f(\mathbf{x}_i; \mathbf{w}) = \mathbf{x}_i \mathbf{w}$
 - Predict \hat{y}_i as $+1$ if $s > 0$, as -1 if $s < 0$
 - If $y_i \neq \hat{y}_i$
 - $\mathbf{w}^{new} = \mathbf{w}^{old} + y \mathbf{x}_i^T$, or
$$\begin{cases} w_0^{new} = w_0^{old} + y x_{i,0} \\ w_1^{new} = w_1^{old} + y x_{i,1} \\ w_2^{new} = w_2^{old} + y x_{i,2} \end{cases}$$

Perceptron update rule (2-D points)

- Iterate over all the points

- Current point is \mathbf{x}_i , and its label is y_i

- Call decision function $s = f(\mathbf{x}_i; \mathbf{w}) = \mathbf{x}_i \mathbf{w}$

- Predict \hat{y}_i as $+1$ if $s > 0$, as -1 if $s < 0$

- If $y_i \neq \hat{y}_i \Leftrightarrow y_i, \hat{y}_i = -1, +1$ or $+1, -1 \Leftrightarrow y_i \hat{y}_i < 0 \Leftrightarrow y_i s < 0$

- $\mathbf{w}^{new} = \mathbf{w}^{old} + y \mathbf{x}_i^T$, or
$$\begin{cases} w_0^{new} = w_0^{old} + y x_{i,0} \\ w_1^{new} = w_1^{old} + y x_{i,1} \\ w_2^{new} = w_2^{old} + y x_{i,2} \end{cases}$$

Perceptron update rule (2-D points)

- Iterate over all the points
 - Current point is \mathbf{x}_i , and its label is y_i
 - If misclassified $\Leftrightarrow yf(\mathbf{x}_i; \mathbf{w}) = y\mathbf{x}_i\mathbf{w} < 0$
 - $\mathbf{w}^{new} = \mathbf{w}^{old} + y\mathbf{x}_i^T$, or
$$\begin{cases} w_0^{new} = w_0^{old} + yx_{i,0} \\ w_1^{new} = w_1^{old} + yx_{i,1} \\ w_2^{new} = w_2^{old} + yx_{i,2} \end{cases}$$
- The update rule for implementation
- We can further define a loss function, but only for theoretic analysis.

Define a loss function

- For misclassified instance \mathbf{x} , $\mathbf{w}^{new} = \mathbf{w}^{old} + y\mathbf{x}^T$ so that
 - $s = f(\mathbf{x}; \mathbf{w}) = \mathbf{x}\mathbf{w}$ is increased if $y = +1$
 - $s = f(\mathbf{x}; \mathbf{w}) = \mathbf{x}\mathbf{w}$ is decreased if $y = -1$
- Loss functions should be minimized
- Define $L(f(\mathbf{x}; \mathbf{w}), y) = -yf(\mathbf{x}; \mathbf{w})$ for misclassified \mathbf{x}
 - $-yf(\mathbf{x}; \mathbf{w}) = -\mathbf{x}\mathbf{w}$ is decreased if $y = +1$
 - $-yf(\mathbf{x}; \mathbf{w}) = +\mathbf{x}\mathbf{w}$ is decreased if $y = -1$
- So $L(f(\mathbf{x}; \mathbf{w}), y) = -yf(\mathbf{x}; \mathbf{w})$ decreases after updating \mathbf{w}

Gradient of the loss function

$$\square L(f(\mathbf{x}; \mathbf{w}), y) = -yf(\mathbf{x}; \mathbf{w}) = -y(x_0w_0 + x_1w_1 + x_2w_2)$$

$$\frac{\partial L}{\partial w_0} = -yx_0 \quad \frac{\partial L}{\partial w_1} = -yx_1 \quad \frac{\partial L}{\partial w_2} = -yx_2$$

\square So

$$\frac{\partial L}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix} = -y \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = -y\mathbf{x}^T$$

The uniform update rule

□ For misclassified instance \mathbf{x} , $\mathbf{w}^{new} = \mathbf{w}^{old} + y\mathbf{x}^T$

□ For misclassified instance \mathbf{x} , the loss function is defined as $L(f(\mathbf{x}; \mathbf{w}), y) = -yf(\mathbf{x}; \mathbf{w})$

$$\frac{\partial L}{\partial \mathbf{w}} = -y\mathbf{x}^T$$

□ So the update rule can be written as

$$\mathbf{w}^{new} = \mathbf{w}^{old} + y\mathbf{x}^T = \mathbf{w}^{old} - \frac{\partial L}{\partial \mathbf{w}}$$

Perceptron update rule (2-D points)

- Define $L(f(\mathbf{x}; \mathbf{w}), y) = -yf(\mathbf{x}; \mathbf{w})$ for misclassified instances
- Iterate over all the points
 - Current point is \mathbf{x}_i , and its label is y_i
 - If misclassified $\Leftrightarrow yf(\mathbf{x}_i; \mathbf{w}) = y\mathbf{x}_i\mathbf{w} < 0 \Leftrightarrow L(f(\mathbf{x}; \mathbf{w}), y) > 0$

$$\mathbf{w}^{new} = \mathbf{w}^{old} - \frac{\partial L}{\partial \mathbf{w}}$$

Perceptron \rightarrow stochastic gradient descent

□ Define

$$L(f(\mathbf{x}; \mathbf{w}), y) = \begin{cases} -yf(\mathbf{x}; \mathbf{w}) & \text{if } L(f(\mathbf{x}; \mathbf{w}), y) > 0 \Leftrightarrow \text{misclassified} \\ 0 & \text{otherwise} \end{cases}$$

□ Then

$$\frac{\partial L}{\partial \mathbf{w}} = \begin{cases} -y\mathbf{x}^T & \text{if } L(f(\mathbf{x}; \mathbf{w}), y) > 0 \Leftrightarrow \text{misclassified} \\ \mathbf{0} & \text{otherwise} \end{cases}$$

□ Iterate over all the points

□ Current point is \mathbf{x}_i , and its label is y_i

□ Apply the update rule

$$\mathbf{w}^{new} = \mathbf{w}^{old} - \frac{\partial L}{\partial \mathbf{w}}$$

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Supervised learning as an optimization problem
 - ❑ Perceptron update rule & loss function
 - ❑ **Logistic regression**
 - ❑ Predict function
 - ❑ Log loss (a.k.a. cross entropy)
- ❑ Notebook tasks
 - ❑ Task 1: Logistic regression
 - ❑ Task 2: Perceptron classifier

Logistic regression (2-D points, 3 classes)

□ Data points

$$\square \mathbf{x}_1 = [1 \quad x_{1,1} \quad x_{1,2}] \rightarrow y_1 = 0, \mathbf{x}_2 = [1 \quad x_{2,1} \quad x_{2,2}] \rightarrow y_2 = 1$$

$$\square \mathbf{x}_3 = [1 \quad x_{3,1} \quad x_{3,2}] \rightarrow y_3 = 2, \mathbf{x}_4 = [1 \quad x_{4,1} \quad x_{4,2}] \rightarrow y_4 = 2$$

□ Model parameters

$$\square \mathbf{W} = [\mathbf{w}_0 \quad \mathbf{w}_1 \quad \mathbf{w}_2] = \begin{bmatrix} w_{0,0} & w_{0,1} & w_{0,2} \\ w_{1,0} & w_{1,1} & w_{1,2} \\ w_{2,0} & w_{2,1} & w_{2,2} \end{bmatrix}$$

□ Decision function, predict \mathbf{x} as class j if s_j is the largest.

$$\square f(\mathbf{x}; \mathbf{w}) = \mathbf{xW} = \begin{bmatrix} \mathbf{xw}_0 \\ \mathbf{xw}_1 \\ \mathbf{xw}_2 \end{bmatrix}^T = \begin{bmatrix} x_0 w_{0,0} + x_1 w_{1,0} + x_2 w_{2,0} \\ x_0 w_{0,1} + x_1 w_{1,1} + x_2 w_{2,1} \\ x_0 w_{0,2} + x_1 w_{1,2} + x_2 w_{2,2} \end{bmatrix}^T = \begin{bmatrix} s_0 \\ s_1 \\ s_2 \end{bmatrix}^T$$

Logistic regression (2-D points, 3 classes)

□ Model parameters

$$\square \mathbf{W} = [\mathbf{w}_0 \quad \mathbf{w}_1 \quad \mathbf{w}_2] = \begin{bmatrix} w_{0,0} & w_{0,1} & w_{0,2} \\ w_{1,0} & w_{1,1} & w_{1,2} \\ w_{2,0} & w_{2,1} & w_{2,2} \end{bmatrix}$$

□ Decision function, predict \mathbf{x} as class j if s_j is the largest.

$$\square f(\mathbf{x}; \mathbf{W}) = \mathbf{xW} = \begin{bmatrix} \mathbf{xw}_0 \\ \mathbf{xw}_1 \\ \mathbf{xw}_2 \end{bmatrix}^T = \begin{bmatrix} x_0 w_{0,0} + x_1 w_{1,0} + x_2 w_{2,0} \\ x_0 w_{0,1} + x_1 w_{1,1} + x_2 w_{2,1} \\ x_0 w_{0,2} + x_1 w_{1,2} + x_2 w_{2,2} \end{bmatrix}^T = \begin{bmatrix} s_0 \\ s_1 \\ s_2 \end{bmatrix}^T$$

□ Output the distribution

$$\square p(y|\mathbf{x}; \mathbf{W}) = \frac{1}{e^{s_0} + e^{s_1} + e^{s_2}} [e^{s_0} \quad e^{s_1} \quad e^{s_2}] \propto [e^{s_0} \quad e^{s_1} \quad e^{s_2}]$$

Logistic regression (2-D points, 3 classes)

□ Decision function, predict \mathbf{x} as class j if s_j is the largest.

$$\square f(\mathbf{x}; \mathbf{W}) = \mathbf{xW} = \begin{bmatrix} \mathbf{xw}_0 \\ \mathbf{xw}_1 \\ \mathbf{xw}_2 \end{bmatrix}^T = \begin{bmatrix} x_0w_{0,0} + x_1w_{1,0} + x_2w_{2,0} \\ x_0w_{0,1} + x_1w_{1,1} + x_2w_{2,1} \\ x_0w_{0,2} + x_1w_{1,2} + x_2w_{2,2} \end{bmatrix}^T = \begin{bmatrix} s_0 \\ s_1 \\ s_2 \end{bmatrix}^T$$

□ Output the distribution

$$\square p(y|\mathbf{x}; \mathbf{W}) = \frac{1}{e^{s_0} + e^{s_1} + e^{s_2}} [e^{s_0} \quad e^{s_1} \quad e^{s_2}] \propto [e^{s_0} \quad e^{s_1} \quad e^{s_2}]$$

□ The log-loss

$$\begin{aligned} L(\mathbf{x}_i, y_i; \mathbf{W}) &= -\log p(y = y_i | \mathbf{x} = \mathbf{x}_i; \mathbf{W}) = -\log \frac{e^{s_{y_i}}}{e^{s_0} + e^{s_1} + e^{s_2}} \\ &= -s_{y_i} + \log(e^{s_0} + e^{s_1} + e^{s_2}) \end{aligned}$$

Log-loss for each data point & mean loss

i	y_i	$p(y = j \mathbf{x} = \mathbf{x}_i, \mathbf{W})$			log-loss
		$j = 0$	$j = 1$	$j = 2$	
1	0	0.6	0.3	0.1	
2	1	0.2	0.7	0.1	
3	2	0.5	0.3	0.2	
4	2	0.3	0.3	0.4	

Log-loss for each data point & mean loss

i	y_i	$p(y = j \mathbf{x} = \mathbf{x}_i, \mathbf{W})$			log-loss
		$j = 0$	$j = 1$	$j = 2$	
1	0	0.6	0.3	0.1	$-\log 0.6$
2	1	0.2	0.7	0.1	$-\log 0.7$
3	2	0.5	0.3	0.2	$-\log 0.2$
4	2	0.3	0.3	0.4	$-\log 0.4$

□ Mean loss (without regularization)

$$\begin{aligned} L &= \frac{1}{4} \sum_{i=1}^4 -\log p(y = y_i | \mathbf{x} = \mathbf{x}_i; \mathbf{W}) \\ &= -\frac{1}{4} (\log 0.6 + \log 0.7 + \log 0.2 + \log 0.4) \end{aligned}$$

Log-loss for each data point & mean loss

i	y_i	$p(y = j \mathbf{x} = \mathbf{x}_i, \mathbf{W})$			↓ log-loss
		$j = 0$	$j = 1$	$j = 2$	
1	0	0.6 ↑	0.3	0.1	$-\log 0.6$
2	1	0.2	0.7 ↑	0.1	$-\log 0.7$
3	2	0.5	0.3	0.2 ↑	$-\log 0.2$
4	2	0.3	0.3	0.4 ↑	$-\log 0.4$

□ Mean loss (without regularization)

$$\begin{aligned} \downarrow L &= \frac{1}{4} \sum_{i=1}^4 -\log p(y = y_i | \mathbf{x} = \mathbf{x}_i; \mathbf{W}) \\ &= -\frac{1}{4} (\log 0.6 + \log 0.7 + \log 0.2 + \log 0.4) \end{aligned}$$

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Supervised learning as an optimization problem
 - ❑ Perceptron update rule & loss function
 - ❑ Logistic regression
 - ❑ Predict function
 - ❑ Log loss (a.k.a. cross entropy)
- ❑ Notebook tasks
 - ❑ Task 1: Logistic regression
 - ❑ Task 2: Perceptron classifier