### CS 453X: Class 9

Jacob Whitehill

# Multi-class versus binary classification

# Multi-class versus binary classification

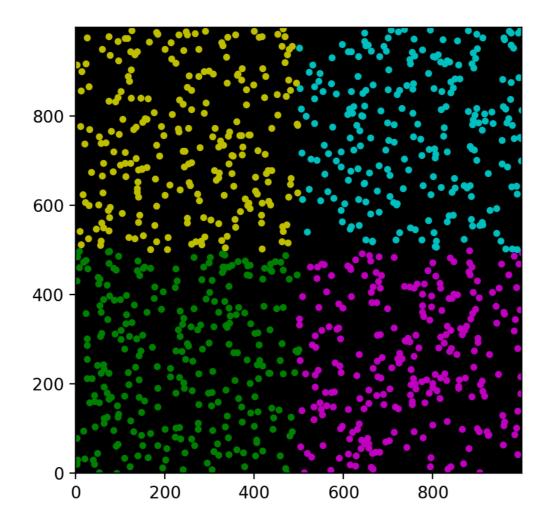
- For the MNIST classification problem, suppose we only really care about distinguishing {0} from {1, 2, ..., 9}.
  - This is essentially a binary classification problem (where {1, 2, ..., 9} become a single "class").
  - But we could still train a softmax regressor on 10 classes.
- Which approach works better?

# Multi-class versus binary classification

- The answer to this question is ultimately empirical there are no theoretical guarantees.
- Intuitive arguments:
  - Binary classification is simpler, exactly matches the problem objective, and requires fewer parameters.
  - Multi-class classification has more parameters and can capture more structure in the data.

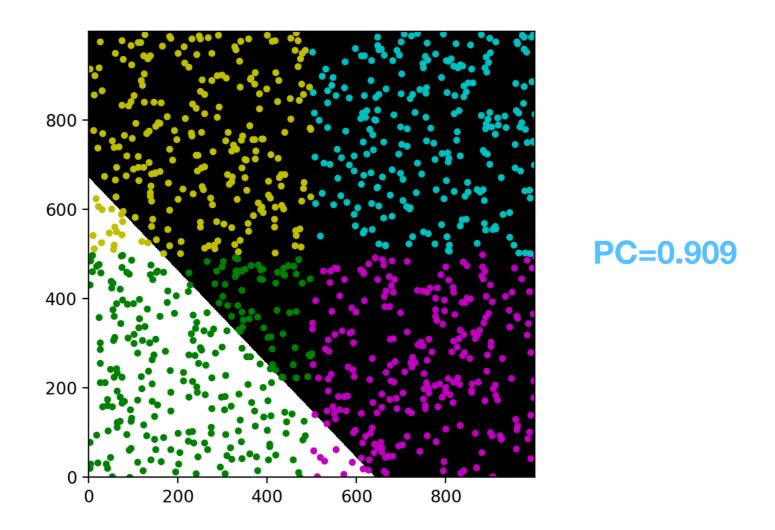
# Example on 4 classes

- Here's a simulation of how multi-class classification can work better.
- Data points are assigned to classes based on the quadrant in the graph.



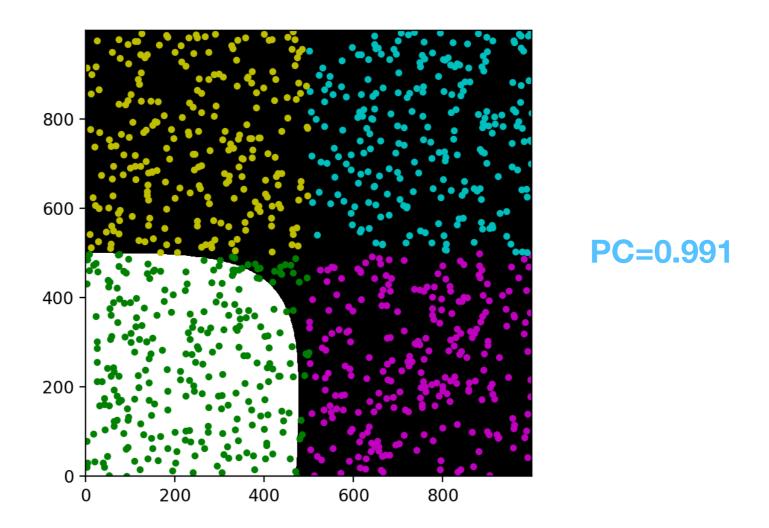
# Example on 4 classes

- Logistic regression (1 v {2,3,4}):
  - White background: set of pixels such that  $\hat{y}_1 > 0.5$



# Example on 4 classes

- Softmax regression (1 v 2 v 3 v 4):
  - White background: set of pixels such that  $\hat{y}_1 > 0.5$



#### Gradient descent

- With gradient descent, we only update the weights after scanning the entire training set.
  - This is slow.
- If the training set contains 60K examples (like in MNIST), then the gradient is an average over 60K images.
  - How much would the gradient really change if we just used, say, 30K images? 15K images? 128 images?

$$\nabla_{\mathbf{W}} f_{\text{CE}}(\mathbf{Y}, \hat{\mathbf{Y}}; \mathbf{W}) = \frac{1}{n} \mathbf{X} (\hat{\mathbf{Y}} - \mathbf{Y})$$

Average over entire training set.

- This is the idea behind stochastic gradient descent (SGD):
  - Randomly sample a small (≪ n) mini-batch (or sometimes just batch) of training examples.
  - Estimate the gradient on just the mini-batch.
  - Update weights based on mini-batch gradient estimate.
  - Repeat.

- In practice, SGD is usually conducted over multiple epochs.
  - An epoch is a single pass through the entire training set.

- In practice, SGD is usually conducted over multiple epochs.
  - An epoch is a single pass through the entire training set.
- Procedure:
  - 1. Let  $\tilde{n} \ll n$  equal the size of the mini-batch.
  - 2. Randomize the order of the examples in the training set.

- In practice, SGD is usually conducted over multiple epochs.
  - An epoch is a single pass through the entire training set.

#### • Procedure:

- 1. Let  $\tilde{n} \ll n$  equal the size of the mini-batch.
- 2. Randomize the order of the examples in the training set.
- 3. For i = 0 to  $(\lceil n/\tilde{n} \rceil 1)$  (one epoch):

- In practice, SGD is usually conducted over multiple epochs.
  - An epoch is a single pass through the entire training set.

#### • Procedure:

- 1. Let  $\tilde{n} \ll n$  equal the size of the mini-batch.
- 2. Randomize the order of the examples in the training set.
- 3. For i = 0 to  $(\lceil n/\tilde{n} \rceil 1)$  (one epoch):
  - A. Select a mini-batch  $\mathcal J$  containing the next  $\tilde n$  examples.

- In practice, SGD is usually conducted over multiple epochs.
  - An epoch is a single pass through the entire training set.

#### Procedure:

- 1. Let  $\tilde{n} \ll n$  equal the size of the mini-batch.
- 2. Randomize the order of the examples in the training set.
- 3. For i = 0 to  $(\lceil n/\tilde{n} \rceil 1)$  (one epoch):
  - A. Select a mini-batch  $\mathcal J$  containing the next  $\tilde n$  examples.
  - B. Compute the gradient on this mini-batch:  $\frac{1}{\tilde{n}}\sum_{i\in\mathcal{J}}\nabla\mathbf{w}f(\mathbf{y}^{(i)},\hat{\mathbf{y}}^{(i)};\mathbf{W})$

- In practice, SGD is usually conducted over multiple epochs.
  - An epoch is a single pass through the entire training set.

#### Procedure:

- 1. Let  $\tilde{n} \ll n$  equal the size of the mini-batch.
- 2. Randomize the order of the examples in the training set.
- 3. For i = 0 to  $(\lceil n/\tilde{n} \rceil 1)$  (one epoch):
  - A. Select a mini-batch  $\mathcal J$  containing the next  $\tilde n$  examples.
  - B. Compute the gradient on this mini-batch:  $\frac{1}{\tilde{n}}\sum_{i\in\mathcal{J}}\nabla\mathbf{w}f(\mathbf{y}^{(i)},\hat{\mathbf{y}}^{(i)};\mathbf{W})$
  - C. Update the weights based on the current mini-batch gradient.

- In practice, SGD is usually conducted over multiple epochs.
  - An epoch is a single pass through the entire training set.

#### Procedure:

- 1. Let  $\tilde{n} \ll n$  equal the size of the mini-batch.
- 2. Randomize the order of the examples in the training set.
- 3. For i = 0 to  $(\lceil n/\tilde{n} \rceil 1)$  (one epoch):
  - A. Select a mini-batch  $\mathcal J$  containing the next  $\tilde n$  examples.
  - B. Compute the gradient on this mini-batch:  $\frac{1}{\tilde{n}}\sum_{i\in\mathcal{J}}\nabla\mathbf{w}f(\mathbf{y}^{(i)},\hat{\mathbf{y}}^{(i)};\mathbf{W})$
  - C. Update the weights based on the current mini-batch gradient.
- 4. Repeat (3) until the desired number of epochs is reached.

- Suppose our training set contains n=8 examples.
- Here is how regular gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.

1
2
3
4
5
6
7
8

- Suppose our training set contains n=8 examples.
- Here is how regular gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - For each round:
    - Compute gradient on all n examples.

1
2
3
4
5
6
7
8

- Suppose our training set contains n=8 examples.
- Here is how regular gradient descent would proceed:
  - Initialize weights **w**<sup>(0)</sup> to random values.
  - For each round:
    - Compute gradient on all n examples.
    - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \boldsymbol{\epsilon} \nabla_{\mathbf{w}} f$

1	
2	
3	
4	
5	
6	
7	
8	

- Suppose our training set contains n=8 examples.
- Here is how regular gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - For each round:
    - Compute gradient on all n examples.
    - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \epsilon \nabla_{\mathbf{w}} f$

1	
2	
3	
4	
5	
6	
7	
8	

- Suppose our training set contains n=8 examples.
- Here is how regular gradient descent would proceed:
  - Initialize weights **w**<sup>(0)</sup> to random values.
  - For each round:
    - Compute gradient on all n examples.
    - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \boldsymbol{\epsilon} \nabla_{\mathbf{w}} f$

1	
2	
3	
4	
5	
6	
7	
8	

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.

1
2
3
4
5
6
7
8

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - Randomize the order of the training data.

4
1
3
5
7
6
8
2

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., *E*): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.

4
1
3
5
7
6
8
2

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., *E*): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \boldsymbol{\varepsilon} \widetilde{\nabla}_{\mathbf{w}} f$

4	
1	
3	
5	
7	
6	
8	
2	

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights **w**<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., *E*): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \epsilon \widetilde{\nabla}_{\mathbf{w}} f$

4
1
3
5
7
6
8
2

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., *E*): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \boldsymbol{\varepsilon} \widetilde{\nabla}_{\mathbf{w}} f$

4	
1	
3	
5	
7	
6	
8	
2	

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights **w**<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., E): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \epsilon \widetilde{\nabla}_{\mathbf{w}} f$

4
1
3
5
7
6
8
2

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., *E*): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \boldsymbol{\varepsilon} \widetilde{\nabla}_{\mathbf{w}} f$

4	
1	
3	
5	
7	
6	
8	
2	

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights **w**<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., E): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \epsilon \widetilde{\nabla}_{\mathbf{w}} f$

4
1
3
5
7
6
8
2

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., *E*): e=1
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \boldsymbol{\varepsilon} \widetilde{\nabla}_{\mathbf{w}} f$

4	
1	
3	
5	
7	
6	
8	
2	

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights **w**<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., *E*): e=2
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \epsilon \widetilde{\nabla}_{\mathbf{w}} f$

4
1
3
5
7
6
8
2

- Suppose our training set contains n=8 examples with  $\tilde{n}=2$ .
- Here is how stochastic gradient descent would proceed:
  - Initialize weights w<sup>(0)</sup> to random values.
  - Randomize the order of the training data.
  - For each epoch (e=1, ..., E): e=2
    - For each round ( $r=1, ..., \lceil n/\tilde{n} \rceil$ ):
      - Compute gradient on next  $\tilde{n}$  examples.
      - Update weights:  $\mathbf{w}^{(t+1)} \longleftarrow \mathbf{w}^{(t)} \epsilon \widetilde{\nabla}_{\mathbf{w}} f$

#### Training examples

4
1
3
5
7
6
8
2

---

- Despite "noise" (statistical inaccuracy) in the mini-batch gradient estimates, we will still converge to local minimum.
- Training can be much faster than regular gradient descent because we adjust the weights many times per epoch.

# SGD: learning rates

- With SGD, our learning rate  $\epsilon$  needs to be **annealed** (reduced slowly over time) to guarantee convergence.
  - Otherwise we might just oscillate forever in weight space.
- Necessary conditions:

$$\lim_{T \to \infty} \sum_{t=1}^{T} |\epsilon_t|^2 < \infty$$

Not too big: sum of squared learning rates converges.

# SGD: learning rates

- With SGD, our learning rate  $\epsilon$  needs to be **annealed** (reduced slowly over time) to guarantee convergence.
  - Otherwise we might just oscillate forever in weight space.
- Necessary conditions:

$$\lim_{T \to \infty} \sum_{t=1}^{T} |\epsilon_t|^2 < \infty \qquad \qquad \lim_{T \to \infty} \sum_{t=1}^{T} |\epsilon_t| = \infty$$

Not too small: sum of absolute learning rates grows to infinity.

# SGD: learning rates

- One common learning rate "schedule" is to multiply  $\epsilon$  by  $c \in (0,1)$  every k rounds.
  - This is called exponential decay.
- Another possibility (which avoids the issue) is to set the number of epochs T to a finite number.
  - SGD may not fully converge, but the machine might still perform well.
- There are many other strategies.

## Convex ML models

## Convex ML models

 The two main ML models we have examined — linear regression and softmax regression — have loss functions that are convex.

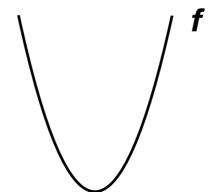
 With a convex function f, every local minimum is also a global minimum.



Convex functions are ideal for conducting gradient descent.

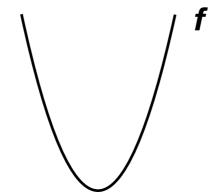
non-convex

How can we tell if a 1-d function f is convex?



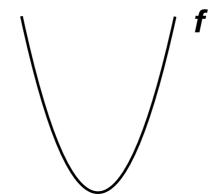
 What property of f ensures there is only one local minimum?

How can we tell if a 1-d function f is convex?



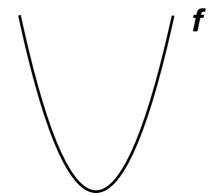
- What property of f ensures there is only one local minimum?
  - From left to right, the slope of f never decreases.

How can we tell if a 1-d function f is convex?



- What property of f ensures there is only one local minimum?
  - From left to right, the slope of f never decreases.
    => the derivative of the slope is always non-negative.

How can we tell if a 1-d function f is convex?



- What property of f ensures there is only one local minimum?
  - From left to right, the slope of *f never decreases*.
    - ==> the derivative of the slope is always non-negative.
    - ==> the second derivative of f is always non-negative.

# Convexity in higher dimensions

• For higher-dimensional *f*, convexity is determined by the second derivative matrix, known as the **Hessian** of *f*.

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \, \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \, \partial x_n} \\ \\ \frac{\partial^2 f}{\partial x_2 \, \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \, \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \\ \frac{\partial^2 f}{\partial x_n \, \partial x_1} & \frac{\partial^2 f}{\partial x_n \, \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}.$$

• For  $f: \mathbb{R}^m \to \mathbb{R}$ , f is convex if the Hessian matrix is positive semi-definite for *every* input **x**.

#### Positive semi-definite

- Positive semi-definite is the matrix analog of being "nonnegative".
- A real symmetric matrix A is positive semi-definite (PSD) if (equivalent conditions):

#### Positive semi-definite

- Positive semi-definite is the matrix analog of being "nonnegative".
- A real symmetric matrix A is positive semi-definite (PSD) if (equivalent conditions):
  - All its eigenvalues are ≥0.
    - If A happens to be diagonal, then its eigenvalues are the diagonal elements.

#### Positive semi-definite

- Positive semi-definite is the matrix analog of being "nonnegative".
- A real symmetric matrix A is positive semi-definite (PSD) if (equivalent conditions):
  - All its eigenvalues are ≥0.
    - If A happens to be diagonal, then its eigenvalues are the diagonal elements.
  - For every vector  $\mathbf{v}$ :  $\mathbf{v}^{\mathsf{T}}\mathbf{A}\mathbf{v} \geq 0$ 
    - Therefore: If there exists any vector v such that
      v<sup>T</sup>Av < 0, then A is not PSD.</li>

## Example

- Suppose  $f(x, y) = 3x^2 + 2y^2 2$ .
- Then the first derivatives are:  $\frac{\partial f}{\partial x} = 6x$   $\frac{\partial f}{\partial y} = 4y$
- The Hessian matrix is therefore:

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x \partial x} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial u \partial x} & \frac{\partial^2 f}{\partial u \partial y} \end{bmatrix} = \begin{bmatrix} 6 & 0 \\ 0 & 4 \end{bmatrix}$$

- Notice that **H** for this f does not depend on (x,y).
- Also, H is a diagonal matrix (with 6 and 4 on the diagonal).
  Hence, the eigenvalues are just 6 and 4. Since they are both non-negative, then f is convex.

## Example

• Graph of  $f(x, y) = 3x^2 + 2y^2 - 2$ :

