# Optimizers  *← Batch size 1*

## Stochastic Gradient Descent (SGD)

Stochastic gradient descent (SGD) in contrast performs a parameter update for *each* training example $x^{(i)}$ and label $y^{(i)}$:

$$\theta = \theta - \eta \cdot \nabla_\theta J(\theta; x^{(i)}; y^{(i)}).$$

Batch gradient descent performs redundant computations for large datasets, as it recomputes gradients for similar examples before each parameter update. SGD does away with this redundancy by performing one update at a time. It is therefore usually much faster and can also be used to learn online. SGD performs frequent updates with a high variance that cause the objective function to fluctuate heavily as in Image 1.
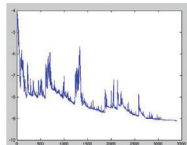


Image 1: SGD fluctuation (Source: Wikipedia)

While batch gradient descent converges to the minimum of the basin the parameters are placed in, SGD's fluctuation, on the one hand, enables it to jump to new and potentially better local minima. On the other hand, this ultimately complicates convergence to the exact minimum, as SGD will keep overshooting. However, it has been shown that when we slowly decrease the learning rate, SGD shows the same convergence behaviour as batch gradient descent, almost certainly converging to a local or the global minimum for non-convex and convex optimization respectively.

## Adaptive Moment (Adam)

*Adam*,[17] which stands for *adaptive moment estimation*, combines the ideas of momentum optimization and RMSProp: just like momentum optimization, it keeps track of an exponentially decaying average of past gradients; and just like RMSProp, it keeps track of an exponentially decaying average of past squared gradients (see Equation 11-8).[18]

*Equation 11-8. Adam algorithm*

1. $\mathbf{m} \leftarrow \beta_1 \mathbf{m} - (1 - \beta_1) \nabla_\theta J(\theta)$
2. $\mathbf{s} \leftarrow \beta_2 \mathbf{s} + (1 - \beta_2) \nabla_\theta J(\theta) \otimes \nabla_\theta J(\theta)$
3. $\widehat{\mathbf{m}} \leftarrow \dfrac{\mathbf{m}}{1 - \beta_1^{\,t}}$
4. $\widehat{\mathbf{s}} \leftarrow \dfrac{\mathbf{s}}{1 - \beta_2^{\,t}}$
5. $\theta \leftarrow \theta + \eta\, \widehat{\mathbf{m}} \oslash \sqrt{\widehat{\mathbf{s}} + \varepsilon}$

*(handwritten: $0 < \beta_1 < 1$, $0 < \beta_2 < 1$; element wise multiplication; time step t; element wise division)*

[17] Diederik P. Kingma and Jimmy Ba, "Adam: A Method for Stochastic Optimization," arXiv preprint arXiv: 1412.6980 (2014).

[18] These are estimations of the mean and (uncentered) variance of the gradients. The mean is often called the *first moment* while the variance is often called the *second moment*, hence the name of the algorithm.

*source*: Géron, Hands-on ML

Others: momentum, RMSprop, Adadelta, Adamax etc.

# Hyperparameters

- A **hyperparameter** is a parameter whose value is used to control the learning process (it is diff**eren**t from "**weights**" or "**parameters**")

- e.g. the batch size is a hyperparameter of gradient descent that controls the number of training samples to work through before the model's internal parameters or weights are updated

- e.g. the number of epochs is a hyperparameter that controls the number of complete passes through the training dataset

- Other hyperparameters: optimizer (e.g. SGD, Adam) and its parameter settings, no. of layers and neurons in each layer, activation functions etc.

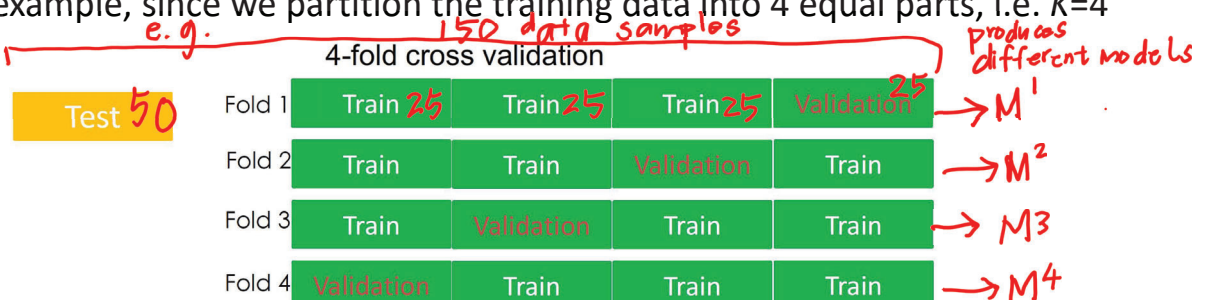# Training, Validation and Test sets

- **Basic method**
- Split entire dataset into **training set** and **validation sets**
  (we may also wish to set aside a **test set**)
  - The training set is used for training
  - The validation set (fixed) is to **simulate the unseen test data** for **model tuning** during the **training stage**
- Use validation set to decide when to stop training and prevent overfitting
- also used for hyperparameter tuning (see later)
- *Strictly speaking, the **test set** is untouched until the final evaluation.*
  *note: sometimes, the "test set" is used as the validation set*
- Performance metrics such as mean-squared error, classification accuracy etc.

---

# Training, Validation and Test sets

- **Advanced method with *K*-fold Cross Validation**
- Split entire dataset into **training set** and **test set**
- The **training set** is used for training and validation
- *The **test set** is untouched until the final evaluation*
- Partition the training set and select one partition to be the validation set
  - next, the validation set is rotated until all the training data have been used for validation
  - in this example, since we partition the training data into 4 equal parts, i.e. *K*=4

e.g.   150 data samples                    produces different models

4-fold cross validation

| Test 50 | | | | |
|---------|---|---|---|---|
| Fold 1 | Train 25 | Train 25 | Train 25 | Validation 25 → M¹ |
| Fold 2 | Train | Train | Validation | Train → M² |
| Fold 3 | Train | Validation | Train | Train → M3 |
| Fold 4 | Validation | Train | Train | Train → M⁴ |

# Hyperparameter Optimization / Tuning

- The validation set is used for tuning the hyperparameters during the training stage
  - evaluate the model with different hyperparameters on the validation set

**Which hyperparameter value(s) should we use?**

*e.g.*
→ 2 hidden nodes
5 hidden nodes

| | Acc. on Val Set 1  Fold 1 | Acc. on Val Set 2  Fold 2 | Acc. on Val Set 3  Fold 3 | Acc. on Val Set 4  Fold 4 | Avg. Acc on Val. Set |
|---|---|---|---|---|---|
| Classifier with hyper-Param1 | 88%  $M_1^1$ | 89%  $M_1^2$ | 93%  $M_1^3$ | 92%  $M_1^4$ | 90.5%  ✔ |
| Classifier with hyper-Param2 | 90%  $M_2^1$ | 88%  $M_2^2$ | 91%  $M_2^3$ | 91%  $M_2^4$ | 90% |

- Hence, select <u>hyperparameter 1</u> and retrain using the entire training set to obtain the final model, which can then be evaluated on the test set
- Note: Hyperparameter optimization will be covered in greater detail later

---

# Performance Metrics: Regression

Evaluation is usually done with validation or test samples, but sometimes we evaluate for training samples

**Mean Square Error**

$$(MSE = \frac{\Sigma_{i=1}^{n}(y_i - \hat{y}_i)^2}{n})$$

**Mean Absolute Error**

$$(MAE = \frac{\Sigma_{i=1}^{n}|y_i - \hat{y}_i|}{n})$$

where $y_i$ denotes the target output and $\hat{y}_i$ denotes the predicted output for sample $i$.



Refer to sklearn.metrics and Keras metrics

# Performance Metrics: Binary Classification

## Confusion Matrix for Binary Classification

Evaluation is usually done with validation or test samples, but sometimes we evaluate for training samples

|  | $\hat{P}$ (predicted) | $\hat{N}$ (predicted) |  |
|---|---|---|---|
| **P (actual)** 20 | True Positive TP 18 | False Negative FN 2 | Sensitivity/ Recall TP/(TP+FN) |
| **N (actual)** 80 | False Positive FP 5 | True Negative TN 75 | Specificity TN/(TN+FP) |
|  | Precision TP/(TP+FP) | Accuracy (TP+TN)/(TP+TN+FP+FN) |  |

Positive
Negative

Refer to sklearn.metrics and Keras metrics

---

# TensorFlow2 and Keras



Useful links:
https://keras.io
https://www.tensorflow.org
https://www.tensorflow.org/api_docs/python/tf/keras
https://www.tensorflow.org/guide/keras/sequential_model
https://www.tensorflow.org/tensorboard/get_started

# *Thank You*
# Questions?

Assoc Prof <u>Tham</u> Chen Khong

E-mail: eletck@nus.edu.sg