Linear Models KMI 2020 @ Nagoya University





Kazu (SLAC) on behalf of the core team (me + Aashwin Ananda Mishra + Francois Drielsma)

Original image credit: xkcd

Useful references

- Online textbook: <u>deep learning (I. Goodfellow and Y. Bengio and A. Courville)</u>
- Online textbook: <u>pattern recognition and machine learning (C. M. Bishop)</u>

... sure, you can log out from this session and start reading those books ;) ...

Machine Learning

- Study of an algorithm that is able to *learn* from data.
 - Given some tasks T and performance measure P, the algorithm is said to learn from data, or experience E, if its performance at T, measured by P, improves with E.
- A cross-road of statistics (probability) and computer science (algorithms) where learning is casted to an optimization problem



Supervised learning

- Given data X and label Y + assume an underlying function P(X)=Y, learn an approximate function Q that mimics P.
- Classification (i.e. discrete labels like dog v.s. cat), regression (i.e. continuous variable like energy of a particle), etc.



Supervised learning

- Given data X and label Y + assume an underlying function P(X)=Y, learn an approximate function Q that mimics P.
- Classification (i.e. discrete labels like dog v.s. cat), regression (i.e. continuous variable like energy of a particle), etc.

Unsupervised learning

- Only given data, learn underlying structure
- Partitioning (clustering), principal component analysis, generative models







Image credit: Nina Miolane, Frédéric Poitevin'

Supervised learning

- Given data X and label Y + assume an underlying function P(X)=Y, learn an approximate function Q that mimics P.
- Classification (i.e. discrete labels like dog v.s. cat), regression (i.e. continuous variable like energy of a particle), etc.

Unsupervised learning

- Only given data, learn underlying structure
- Partitioning (clustering), principal component analysis, generative models

Reinforcement learning

• May not even have static data! Learn to gain most cumulative reward by interacting with the environment



Supervised learning

- Given data X and label Y + assume an underlying function P(X)=Y, learn an approximate function Q that mimics P.
- Classification (i.e. discrete labels like dog v.s. cat), regression (i.e. continuous variable like energy of a particle), etc.

Unsupervised learning

- Only given data, learn underlying structure
- Partitioning (clustering), principal component analysis, generative models

Reinforcement learning

- May not even have static data! Learn to gain most cumulative reward by interacting with the environment
- Playing a game, facility control, etc.



Supervised learning

- Given data X and label Y + assume an underlying function P(X)=Y, learn an approximate function Q that mimics P.
- Classification (i.e. discrete labels like dog v.s. cat), regression (i.e. continuous variable like energy of a particle), etc.



Supervised learning

- Given data X and label Y + assume an underlying function P(X)=Y, learn an approximate function Q that mimics P.
- Classification (i.e. discrete labels like dog v.s. cat), regression (i.e. continuous variable like energy of a particle), etc.



Unknown function P produce the response (y₀, y₁, ...) from data (x₀, x₁, ...)



Unknown function P produce the response (y₀, y₁, ...) from data (x₀, x₁, ...)
1. Define a parameterized model Q to mimic the behavior of P.



Unknown function P produce the response (y₀, y₁, ...) from data (x₀, x₁, ...)
Define a parameterized model Q to mimic the behavior of P.

2. Tune the parameters of Q to minimize the divergence between P and Q.



Unknown function P produce the response (y₀, y₁, ...) from data (x₀, x₁, ...)
Define a parameterized model Q to mimic the behavior of P.

2. Tune the parameters of Q to minimize the divergence between P and Q.





Unknown function P produce the response $(y_0, y_1, ...)$ from data $(x_0, x_1, ...)$ 1. Define a parameterized model Q to mimic the behavior of P.

- 2. Tune the parameters of Q to minimize the divergence between P and Q.
 - \circ Define a loss function
 - Minimize the loss





Full disclosure about data (left):

• x sampled from flat distribution [0,10]

•
$$y = a^*x + b + \varepsilon$$

 $\circ~$ a=2, b=1.5, ϵ ~ Normal (µ=0, σ =1.0)

In general, linear regression employs

$$Q(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i}^{N} w_i \psi_i(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \cdot \boldsymbol{\Psi}$$



Full disclosure about data (left):

• x sampled from flat distribution [0,10]

•
$$y = a^*x + b + \varepsilon$$

 $\circ~$ a=2, b=1.5, ϵ ~ Normal (µ=0, σ =1.0)

In general, linear regression employs

$$Q(\mathbf{x}, \mathbf{w}) = w_0 + \sum_i^N w_i \psi_i(\mathbf{x}) = \mathbf{w}^T \cdot \mathbf{\Psi}$$

A simple, intuitive model: $Q(x) = w_0 + w_1 x$



Full disclosure about data (left):

• x sampled from flat distribution [0,10]

•
$$y = a^*x + b + \varepsilon$$

 $\circ~$ a=2, b=1.5, ϵ ~ Normal (µ=0, σ =1.0)

In general, linear regression employs

$$Q(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i}^{N} w_i \psi_i(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \cdot \boldsymbol{\Psi}$$

A simple, intuitive model: $Q(x) = w_0 + w_1 x$

Popular one: Mean Square Error (MSE) $\mathcal{L} = \frac{1}{m} \sum_{i=1}^{m} (y_i - Q(x_i))^2$

Tune **w** for:
$$\nabla_{\mathbf{w}} \mathcal{L} = 0$$



Goal: tune the parameters **w** and achieve $\nabla_{\mathbf{w}} \mathcal{L} = 0$

Goal: tune the parameters **w** and achieve $\nabla_{\mathbf{w}} \mathcal{L} = 0$

Minimize the loss iteratively:
$$\mathbf{w}_{i+1} = \mathbf{w}_i - \lambda
abla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{x})$$

called **Gradient Descent** (GD) where λ controls the rate of learning process.

Goal: tune the parameters **w** and achieve $\nabla_{\mathbf{w}} \mathcal{L} = 0$

Minimize the loss iteratively:
$$\mathbf{w}_{i+1} = \mathbf{w}_i - \lambda
abla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{x})$$

called **Gradient Descent** (GD) where λ controls the rate of learning process.

For a smooth convex functions, GD achieves the loss O(1/k) after k steps. For non-convex functions, there may be multiple optimal points.





Goal: tune the parameters **w** and achieve $\nabla_{\mathbf{w}} \mathcal{L} = 0$... **faster**?

Mini-batch Stochastic GD (SGD) use a subset of data for gradient.

Goal: tune the parameters **w** and achieve $\nabla_{\mathbf{w}} \mathcal{L} = 0$... **faster**?

Mini-batch Stochastic GD (SGD) use a subset of data for gradient.

- 1. Create a batch = random subset of data.
- 2. Compute the gradient for the batch and update the parameters. Note: when data is always new (never seen before), called "online learning"





Goal: tune the parameters **w** and achieve $\nabla_{\mathbf{w}} \mathcal{L} = 0$... **faster**?

Mini-batch Stochastic GD (SGD) use a subset of data for gradient.

- 1. Create a batch = random subset of data.
- 2. Compute the gradient for the batch and update the parameters. Note: when data is always new (never seen before), called "online learning"





Basic loss functions for regressions 1. Mean Absolute Error (MAE, L1 loss) • Could benefit from an MAE = $\frac{1}{m} \sum_{i=1}^{m} |y_i - Q(\mathbf{x}, \mathbf{w})|$ • Could benefit from an adjusted learning rate 2. Mean Square Error (MSE, L2 loss) $MSE = \frac{1}{m} \sum_{i=1}^{m} (y_i - Q(x_i))^2$ • Loss gets small when <1 but may explode when >> 1

3. Huber Loss ... combine them together

Huber =
$$\begin{cases} \frac{1}{2}a^2 \dots \text{for } a \leq \delta\\ \delta|a| - \frac{1}{2}\delta^2 \dots \text{otherwise} \end{cases}$$

L1 while loss is large, L2 when it's small (δ : *hyperparameter*)

(will be back on "hyperparameter" later)

A metric to measure the divergence between P and Q

$$D_{\mathrm{KL}}(P||Q) = \int P(\mathbf{x}, y) \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} d\mathbf{x} dy$$

Kullback-Leibler (KL) Divergence

 $D_{KL} \ge 0$ always, and 0 if and only if P and Q are identical

A metric to measure the divergence between P and Q

$$D_{\mathrm{KL}}(P \| Q) = \int P(\mathbf{x}, y) \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} d\mathbf{x} dy$$

Kullback-Leibler (KL) DivergenceDKL >= 0 always, and 0 if and only ifP and Q are identical

Minimization of D_{KL} using (S)GD can be "tuning of Q_w "

 $\mathbf{w}_{i+1} = \mathbf{w}_i - \lambda \nabla_{\mathbf{w}} D_{\mathrm{KL}}(P \| Q)$

A metric to measure the divergence between P and Q

$$D_{\mathrm{KL}}(P \| Q) = \int P(\mathbf{x}, y) \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} d\mathbf{x} dy$$

Kullback-Leibler (KL) Divergence

 $D_{KL} \ge 0$ always, and 0 if and only if P and Q are identical

Minimization of D_{KL} using (S)GD can be "tuning of Q_w "

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \lambda \nabla_{\mathbf{w}} D_{\mathrm{KL}}(P \| Q) \Longrightarrow \int P(\mathbf{x}, y) \nabla_{\mathbf{w}} \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} d\mathbf{x} dy = \left\langle \nabla_{\mathbf{w}} \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} \right\rangle$$

A metric to measure the divergence between P and Q

$$D_{\mathrm{KL}}(P \| Q) = \int P(\mathbf{x}, y) \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} d\mathbf{x} dy$$

Kullback-Leibler (KL) Divergence

Dĸ∟ >= 0 always, and 0 if and only if P and Q are identical

Minimization of *D*_{KL} using (S)GD can be "tuning of Q_w"

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \lambda \nabla_{\mathbf{w}} D_{\mathrm{KL}}(P \| Q) \Longrightarrow \int P(\mathbf{x}, y) \nabla_{\mathbf{w}} \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} d\mathbf{x} dy = \left\langle \nabla_{\mathbf{w}} \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} \right\rangle$$

Can compute the expectation value using m data instances

$$\left\langle \nabla_{\mathbf{w}} \log \frac{P(\mathbf{x}, y)}{Q_{\mathbf{w}}(\mathbf{x}, y)} \right\rangle = \frac{1}{m} \sum_{\mathbf{x}, y} \nabla_{\mathbf{w}} \log \frac{P(\mathbf{x}_i, y_i)}{Q_{\mathbf{w}}(\mathbf{x}_i, y_i)}$$

... *P* is still an unknown function, but average probability of $P(x_i, y_i)$ might be possible to approximate from many samples

Consider a measurement

- Total number of measurements: *N*
- Observed occurrence N_i for an event type i

We cannot measure underlying probability p_i but can build a model q_i to approximate p_i , and try to find a "good q_i ".

How to define + find a "good q_i ?"

Consider a measurement

- Total number of measurements: *N*
- Observed occurrence N_i for an event type i

We cannot measure underlying probability p_i but can build a model q_i to approximate p_i , and try to find a "good q_i ".

How to define + find a "good q_i ?" Maximum Likelihood Estimation (MLE) • Maybe: "good" if q_i has a high likelihood to yield the observed data

Consider a measurement

- Total number of measurements: *N*
- Observed occurrence N_i for an event type i

We cannot measure underlying probability p_i but can build a model q_i to approximate p_i , and try to find a "good q_i ".

How to define + find a "good q_i ?"

Maximum Likelihood Estimation (MLE)

- Maybe: "good" if q_i has a high likelihood to yield the observed data
- Probability to observe data x: $P(\vec{x}|\vec{q}) = q_0^{N_0} \cdot q_1^{N_1} \cdots q_k^{N_k} \times \frac{N!}{N_0! \cdot N_1! \cdots N_k!}$
 - When N and N_i are large ... $N_i \approx N \cdot p_i$ and $N_i! \approx N_i^{N_i}$

$$\circ \dots$$
 which yields $P(\vec{x}|\vec{q}) \approx \exp\left[-\sum_{i}^{\text{all}} p_i \log \frac{p_i}{q_i}\right]$

Maximizing this likelihood means to take the limit of the exponent becoming zero, which is same as minimizing $D_{\rm KL}$

Consider a measurement

- Total number of measurements: N
- Observed occurrence N_i for an event type i

We cannot measure underlying probability p_i but can build a model q_i to approximate p_i , and try to find a "good q_i ".

How to define + find a "good q_i ?"

Maximum Likelihood Estimation (MLE)

- Maybe: "good" if q_i has a high likelihood to yield the observed data
- Probability to observe data x: $P(\vec{x}|\vec{q}) = q_0^{N_0} \cdot q_1^{N_1} \cdots q_k^{N_k} \times \frac{N!}{N_0! \cdot N_1! \cdots N_k!}$
 - When N and N_i are large ... $N_i \approx N \cdot p_i$ and $N_i! \approx N_i^{N_i}$ This is a crude, partial credit argument (or zero credit if you got a lucky prof.)

• ... which yields
$$P(\vec{x}|\vec{q}) \approx \exp\left[-\sum_{i}^{\text{all}} p_i \log \frac{p_i}{q_i}\right]$$

Maximizing this likelihood means to take the limit of the exponent becoming zero, which is same as minimizing DKL

Recall our data

- x sampled from flat distribution [0,10]
- $y = a^*x + b + \varepsilon$

 $\circ~$ a=2, b=1.5, ϵ ~ Normal (µ=0, σ =1.0)



Recall our data

- x sampled from flat distribution [0,10]
- $y = a^*x + b + \varepsilon$

 $\circ~$ a=2, b=1.5, ϵ ~ Normal (µ=0, σ =1.0)

... so equivalently we can consider that

$$P(y_i|x_i) \propto \exp\left[-\frac{(y_i - (ax_i + b))^2}{2\sigma^2}\right]$$



Recall our data

- x sampled from flat distribution [0,10]
- $y = a^*x + b + \varepsilon$

 $\circ~$ a=2, b=1.5, ϵ ~ Normal (µ=0, σ =1.0)

... so equivalently we can consider that

$$P(y_i|x_i) \propto \exp\left[-\frac{(y_i - (ax_i + b))^2}{2\sigma^2}\right]$$

The likelihood is a product of all data instances

Likelihood
$$\mathcal{L} = \prod_i P(y_i|x_i)$$



Recall our data

- x sampled from flat distribution [0,10]
- $y = a^*x + b + \varepsilon$

 $\circ~$ a=2, b=1.5, ϵ ~ Normal (µ=0, σ =1.0)

... so equivalently we can consider that

$$P(y_i|x_i) \propto \exp\left[-\frac{(y_i - (ax_i + b))^2}{2\sigma^2}\right]$$



The likelihood is a product of all data instances

Likelihood
$$\mathcal{L} = \prod_{i} P(y_i | x_i) \Rightarrow -\log \mathcal{L} \propto \sum_{i} (y_i - (ax_i + b))^2 \dots (MSE)$$

So minimizing MSE was same as maximum likelihood estimate **in this case**.

Nothing special. We can just generalize for m regression targets...

$$MSE = \sum_{i}^{n} \left(y_i - \sum_{j}^{m} w_j x_{ji} \right)^2 = \|\mathbf{y} - \mathbf{X} \cdot \mathbf{w}\|^2$$

Nothing special. We can just generalize for m regression targets...

$$MSE = \sum_{i}^{n} \left(y_{i} - \sum_{j}^{m} w_{j} x_{ji} \right)^{2} = \|\mathbf{y} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

$$\|\mathbf{x} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

Nothing special. We can just generalize for m regression targets...

$$MSE = \sum_{i}^{n} \left(y_{i} - \sum_{j}^{m} w_{j} x_{ji} \right)^{2} = \|\mathbf{y} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

$$\|\mathbf{x} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

$$\|\mathbf{x} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

$$\|\mathbf{x} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

Taking a gradient and solve for minimization:

$$\nabla_{\mathbf{w}} MSE = 2(\mathbf{y} - \mathbf{w} \cdot \mathbf{X})\mathbf{X} \Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Nothing special. We can just generalize for m regression targets...

$$MSE = \sum_{i}^{n} \left(y_{i} - \sum_{j}^{m} w_{j} x_{ji} \right)^{2} = \|\mathbf{y} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

$$\|\mathbf{x} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

$$\|\mathbf{x} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

$$\|\mathbf{x} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

Taking a gradient and solve for minimization:

$$\nabla_{\mathbf{w}} MSE = 2(\mathbf{y} - \mathbf{w} \cdot \mathbf{X})\mathbf{X} \Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Recall "linear" means "linear to features"… My "features" could be *m*-th power of *x* … Maybe fit 15-th degree polynomial to a sin curve?

Nothing special. We can just generalize for m regression targets...

$$MSE = \sum_{i}^{n} \left(y_{i} - \sum_{j}^{m} w_{j} x_{ji} \right)^{2} = \|\mathbf{y} - \mathbf{X} \cdot \mathbf{w}\|^{2}$$

Taking a gradient and solve for minimization:

$$\nabla_{\mathbf{w}} MSE = 2(\mathbf{y} - \mathbf{w} \cdot \mathbf{X}) \mathbf{X} \Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Recall "linear" means "linear to features"… My "features" could be *m*-th power of *x* … Maybe fit 15-th degree polynomial to a sin curve?



From scikit webpage

Hyperparameter: design constants of algorithms (e.g. learning rate)

Overfit: a model has "**memorized data**". Works well on train data but poorly on independent samples. Typical for complex model + low data statistics

Generalization: model works well equally on the train and unseen datasets

Hyperparameter: design constants of algorithms (e.g. learning rate)

Overfit: a model has "**memorized data**". Works well on train data but poorly on independent samples. Typical for complex model + low data statistics

Generalization: model works well equally on the train and unseen datasets



In this case...

A polynomial of a higher power (hyperparameter!) makes the model more complex, or flexible, and as a result a model can overfit.

Cross-Validation

Two datasets

- Train set to optimize your model
- Test set to measure performance

... we want to use the test set only once at the end, never want to tune anything on it!

Cross-Validation

Two datasets

- **Train set** to optimize your model
- **Test set** to measure performance

... we want to use the test set only once at the end, never want to tune anything on it!

Strategy:

- Split the train set into k-folds
- Use *i*-th set as a "**validation set**" to measure the performance of a model trained on the rest (k-1 combined).
- Repeat *k*-times and take the mean as a performance.



Regularization

There may be more than one solution (i.e. local minimum in GD)... prefer a *simpler* solution over complicated ones (=may help generalization).

$$\mathcal{L}_{\text{total}} = \mathcal{L}\left(\mathbf{y}, Q(\mathbf{x}, \mathbf{w})\right) + \alpha R(\mathbf{w})$$
model loss regularization loss

Regularization

There may be more than one solution (i.e. local minimum in GD)... prefer a *simpler* solution over complicated ones (=may help generalization).

$$\mathcal{L}_{\text{total}} = \mathcal{L}\left(\mathbf{y}, Q(\mathbf{x}, \mathbf{w})\right) + \alpha R(\mathbf{w})$$
model loss regularization loss

Basic regularization terms are L1, L2, and L1+L2



Regularization

There may be more than one solution (i.e. local minimum in GD)... prefer a *simpler* solution over complicated ones (=may help generalization).

$$\mathcal{L}_{\text{total}} = \mathcal{L}\left(\mathbf{y}, Q(\mathbf{x}, \mathbf{w})\right) + \alpha R(\mathbf{w})$$
model loss regularization loss

Basic regularization terms are L1, L2, and L1+L2

$$L2 = \|\mathbf{w}\|^2 = \sum_{i} w_i^2$$

 $L1 = \|\mathbf{w}\| = \sum |w_i|$

L2 favors weights with smaller values

$$L1 + L2 = \sum_{i} \left(\beta w_i^2 + |w_i|\right)$$

Can also combine ("elastic net")

L1 favors sparse solution (also called LASSO)

Intermediate summary

- "Learning" is a process of tuning a model Q(x) to make a better approximation of inaccessible, underlying function P(x)
 - Define a model and a loss function, use Gradient Descent (GD) to tune the parameters
 - Stochastic GD (SGD) uses a random subset of train data per parameter update
 - Minimization of *D*KL, maximization of likelihood (MLE)
- Linear regression = linear transformation of features (i.e. linear to model parameters)
 - $\circ~$ Exact solution is available and can be compared to the optimization outcome
 - Popular loss functions = Mean Squared Error (MSE), Mean Absolute Error (MAE) ...
- Train set = dataset used to optimize Q(x)
- Test set = dataset used to benchmark the performance of Q(x)
- Generalization = model performs on the test dataset as well as the train dataset
- Overfit = *memorization* of data, a model performs well on train set but poorly on test set
- Cross validation = splitting the train set to k-folds to create validation set(s) and measure model performance and/or tune hyperparameters
- Regularization = additional constraints on model parameters, can help avoiding overfitting
 L2 prefers smaller weight values, L1 may lead to a sparse solution



Full disclosure about data (left):

- Red data points ~ Normal (μ =(-1,-1), σ =(1,1))
- Blue data points ~ Normal (μ =(1,1), σ =(1,1))



Full disclosure about data (left):

- Red data points ~ Normal (μ =(-1,-1), σ =(1,1))
- Blue data points ~ Normal (μ =(1,1), σ =(1,1))

What should be Q and loss function?

- A straight line to separate two colors
- Classification confidence low on and near the line, more confident away from the line



Full disclosure about data (left):

- Red data points ~ Normal (μ =(-1,-1), σ =(1,1))
- Blue data points ~ Normal (μ =(1,1), σ =(1,1))

What should be Q and loss function?

- A straight line to separate two colors
- Classification confidence low on and near the line, more confident away from the line

$$Q(\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^T \cdot \mathbf{x})}$$

Popular one: logistic (a.k.a. "sigmoid") function

• Satisfies what we want + the output is bounded [0,1] (probabilistic)





Full disclosure about data (left):

- Red data points ~ Normal (μ =(-1,-1), σ =(1,1))
- Blue data points ~ Normal (μ =(1,1), σ =(1,1))

What should be Q and loss function?

- A straight line to separate two colors
- Classification confidence low on and near the line, more confident away from the line

$$Q(\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^T \cdot \mathbf{x})}$$

Popular one: logistic (a.k.a. "sigmoid") function

• Satisfies what we want + the output is bounded [0,1] (probabilistic)

Likelihood =
$$\prod_{i}^{m} (y_i Q(\mathbf{x}_i, \mathbf{w}) + (1 - y_i)(1 - Q(\mathbf{x}_i, \mathbf{w}))$$

y_i = 0 if red, 1 if blue





Full disclosure about data (left):

- Red data points ~ Normal (μ =(-1,-1), σ =(1,1))
- Blue data points ~ Normal (μ =(1,1), σ =(1,1))

What should be Q and loss function?

- A straight line to separate two colors
- Classification confidence low on and near the line, more confident away from the line

$$Q(\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^T \cdot \mathbf{x})}$$

Popular one: logistic (a.k.a. "sigmoid") function

• Satisfies what we want + the output is bounded [0,1] (probabilistic)

$$\mathcal{L} = -\sum_{i}^{m} \left[y_i \log(Q(\mathbf{x}_i, \mathbf{w})) + (1 - y_i) \log(1 - Q(\mathbf{x}_i, \mathbf{w})) \right]$$

Minimization target = *Loss* = negative log likelihood

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{1}{m} \mathbf{x}^T \left(Q(\mathbf{x}, \mathbf{w}) - \mathbf{y} \right)$$

For gradient update!





Generalization to multiple targets (>2 categorization classes)

• Softmax function for Q + MLE yields cross-entropy loss function

$$Q_i(\mathbf{x}, \mathbf{w}_i) = \frac{\exp(-\mathbf{w}_i^T \cdot \mathbf{x})}{\sum_j \exp(-\mathbf{w}_j^T \cdot \mathbf{x})}$$

Q*i* is a probability for being *i*-th class. Need a set of **w** per class!

$$\mathcal{L} = -\sum_{i} y_{i} \log Q_{i} \left(\mathbf{x}_{i}, \mathbf{w}_{i} \right)$$

 $y_i = 1$ only for the correct class

How should we compare the classification performance that depends on an interpretation of score (i.e. Q)? This may be application specific, but a standard procedure exists with useful jargons :)

How should we compare the classification performance that depends on an interpretation of score (i.e. Q)? This may be application specific, but a standard procedure exists with useful jargons :)

	Label P=1	Label P=0
Prediction Q=1	True Positive (TP)	False Positive (FP)
Prediction $Q=0$	False Negative (FN)	True Negative (TN)



How should we compare the classification performance that depends on an interpretation of score (i.e. Q)? This may be application specific, but a standard procedure exists with useful jargons :)

	Label P=1	Label P=0
Prediction Q=1	True Positive (TP)	False Positive (FP)
Prediction Q=0	False Negative (FN)	True Negative (TN)



$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

How should we compare the classification performance that depends on an interpretation of score (i.e. Q)? This may be application specific, but a standard procedure exists with useful jargons :)

	Label P=1	Label P=0
Prediction Q=1	True Positive (TP)	False Positive (FP)
Prediction Q=0	False Negative (FN)	True Negative (TN)



$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

False Positive Rate (FPR) =
$$\frac{FP}{FP + TN}$$

True Positive Rate (TPR) = $\frac{TP}{TP + FN}$

We can create a curve of FPR v.s. TPR by varying the ¹⁰ score threshold. This is **Receiver Operating Characteristic (ROC)** TF curve. The area under this curve may be used as a performance metric.



How should we compare the classification performance that depends on an interpretation of score (i.e. Q)? This may be application specific, but a standard procedure exists with useful jargons :)

	Label P=1	Label P=0
Prediction Q=1	True Positive (TP)	False Positive (FP)
Prediction Q=0	False Negative (FN)	True Negative (TN)

Precision and Recall are also often used metrics. In physics they are sometimes called "purity" and "efficiency" of the prediction!







More dictionary

- Logistic regression is a task to classify an input into predefined set of classes
 - Logistic loss function gives a probabilistic interpretation of a score for binary classification
 - Generalization to multinomial regression using softmax function and cross-entropy loss.
 - $\circ~$ (S)GD can be used to optimize a linear model
- Classification results with labels can be grouped into 4 categories
 - True Positive or TP... (prediction, label) = (1,1)
 - False Positive or FP ... (prediction, label) = (1,0)
 - False Negative or FN ... (prediction, label) = (0,1)
 - True Negative or TN ... (prediction, label) = (0,0)
 - Accuracy = (TP+TN) / ALL
 - False Positive Rate or FPR = FP / (FP+TN)
 - True Positive Rate or TPR = TP / (TP+FN)
 - $\circ \quad \text{Precision} = \text{TP} / (\text{TP} + \text{FP})$
 - Recall = TP / (TP+FN)
- FPR v.s. TPR = Receiver Operating Characteristic (ROC) curve, the area under the curve is larger for a better classifier
- Precision v.s. Recall curve is another metric that is more robust against class imbalance

Google

nagoya food



restaurants

Q All Settings Tools

aichi



nago

>

Collections SafeSearch -

eat

miso	nikomi	udon

nagoya station

Food In Nagoya, Japan ...

trip101.com

halal



Ultimate Nagova Guide For F... favy-jp.com



Top 9 Famous Foods of Nagoya | Japan Info jpninfo.com



Nagova Food: Top 8 Dishes You Have to ... mariaabroad.com



japan

0

Q

aichi prefecture

The Top 5 Must Try Food in Nagoya blog.gaijinpot.com



Nagoya: Home to Japan's most unique ... cnn.com



yamamotoya

Ultimate Nagova Guide For Food Lovers ... favy-jp.com



noodles

spaghetti

miso katsu

The Top 5 Must Try Food in Nagoya blog.gaijinpot.com

Nagova Food - Top 10 dishes for Japan... eatwanderexplore.com

dishes

japanese

centrip-japan.com

cuisine





Nagoya Meshi ("Must Eat" Local Foods in ... japanfoodculture.org



Nagoya Food: 8 Nagoya Meshi Dishes to ... willflyforfood.net



Nagoya: Home to Japan's most unique ... cnn.com



big destination for food in Japan mic.com



Best Restaurants in Nagoya ...



Ultimate Nagoya Guide For Food Lovers...

favy-jp.com

The Top 5 Must Try Food in Nagoya blog.gaijinpot.com



7 Best Nagoya Foods to Try - Trip-N-Tra... trip-n-travel.com



Must-Eat Foods in Nagova... centrip-japan.com



Let's go to eat yummy food to Nagoya. jpnfood.com



Nagoya Meshi ("Must Eat" Local Foods i... japanfoodculture.org



Nagoya Popular Dish, Miso-Stewed U... fooddiversity.today



5 Restaurants Near Na... wow-j.com



The Top 5 Must Try Food in Nagoya blog.gaijinpot.com



Nagoya Food Guide - 20 Local Foods in ... nagoyafoodie.com



























Nagoya: Home to Japan's most unique ... cnn.com



the mark





















Played with linear models today. What if we convolve multiple linear transformations? **Next topic**: Neural Networks!

Scratch Slides

Core idea

- Consider 2 statements ... which one is information rich?
 - A someone taught gorilla a linear algebra
 - A gorilla taught someone a linear algebra
 - $\circ \ \ ...$ a surprising event contains more information
- Information conte $-\log P(A)$
 - Additive: information from event A + B: $-\log P(A) \cdot P(B) = -[\log P(A) + \log P(B))]$
 - e.g.) "I picked 13 of diamond from the stack": $-\log\left(\frac{1}{13} \cdot \frac{1}{4}\right) = \log 52$
- Total expected information content

- 11

• **Entropy:**
$$\mathbf{S} = -\sum_{i}^{\mathrm{an}} P_i \log P_i$$

Core idea

- Consider a measurement
 - $\circ~$ Total number of measurements: N
 - Observed occurrence N_i for an event type *i*

We cannot measure underlying probability p_i but can build a model q_i to approximate p_i , and say the task for machine learning is to learn a "good q_i ".

How to define + find a "good q_i ?"

- Maybe: "good" if q_i has a high likelihood to vield the observed data
- Probability to observe data x: $P(\vec{x}|\vec{q}) = q_0^{N_0} \cdot q_1^{N_1} \cdots q_k^{N_k} \times \frac{N!}{N_0! \cdot N_1! \cdots N_k!}$
 - When N and N_i are large ... $N_i \approx N \cdot p_i$ and $N_i! \approx N_i^{N_i}$

• ... which yields
$$P(\vec{x}|\vec{q}) \approx \exp\left[-\sum_{i}^{\text{all}} p_i \log \frac{p_i}{q_i}\right]$$

Maximize the likelihood = minimize the