자매품으로 Softmax
(multinomial logistic regression)도 많이 씀
- 특히 DL에서는 Softmax를 더!

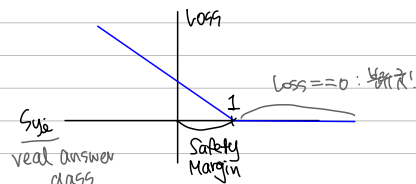○ Multi-class SVM: binary-class SVM의 일반화된 형태
    └두개의 class만 다룸 ↔ CIFAR-10 has 10 classes.
    └ positive/negative

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$ = safety margin.

맞지 않은 카테고리에 대해

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

→ -if) 정답 class의 score 점수 가장 높으면:
then) $\max(0, s_j - s_{y_i} + 1)$
    Max(0, value)의 형태 = "hinge loss"



loss
$s_{y_i}$ / real answer class
1 / Safety Margin
loss==0: 불만족!

Q. Safety Margin?
스코어가 상대적인 차이를 강조함.
서로 너무 가깝지않게 해줌?? 응

손실함수 $L_i$: 얼마나 구리게 예측? 정량화함.
$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$
최종 Loss "L": 각 N개 sample 들의 Loss 평균.
- W 탐색해서 Loss를 minimize하는 W 찾게 될 것임.

Multi-class SVM loss
이것 SVM의 일반화된 형태
SVM (Support vector machine)

---

Q. What happens to loss if car scores change a bit?
SVM loss는 "다른 스코어들과의 차이" 고려함.
약간 1 이내만의 car score가 상대적으로 영향 크지 때문에 "충분 벗어나서 써보면 Margin 위치
→ loss change X

Q. What is the min/max possible loss?
- minimum: 0
  if every class, answer class score is very big
- maximum: ∞
  consider the loss function is hingle loss shape.
  if answer class score : 엄청 작은 neg value.

Q. At initialization W is small so all s~0.
what is the loss?
: N(class) - 1.
all s~0, and our margin is 1. since, we get the loss: C-1
⇒ useful debugging strategy:
Loss가 이상하게 뛸지 검토 가능. (if 트레이닝 첫 시작 시 Loss ≠ C-1, bug)

Q. what if the sum was over all classes?
(including $j = y_i$)
  정답인 class까지 다 더하면: loss += 1
'원래' $j = y_i$인 class (정답 class): loss == 0. ⇒ loss가 0이면 없다는 것이 없다고 해석 가능하기 때문.

Q. What if we used mean instead of sum?
Just rescaling.
상대적으로 비교하므로 의미 없음.

Q. What if we used $\sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$   Squared hinge loss.
(kind of tradeoff
between good and badness in kind of nonlinear way)
→ different. 분리 경우와 다른 없음

목적에 따라, { hinge loss: squared 보다 안 좋음 차이 덜 완화
상황에 맞게 선택 { squared hinge loss: 엄청 안좋은것 → 더 x100 안 좋은거

## Example Code

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

numpy 이용해서 함수작성 가능 →

```
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1)
    margins[y] = 0   정확한 class만 0으로 만들어줌 (vectorized)
    loss_i = np.sum(margins)
    return loss_i
```
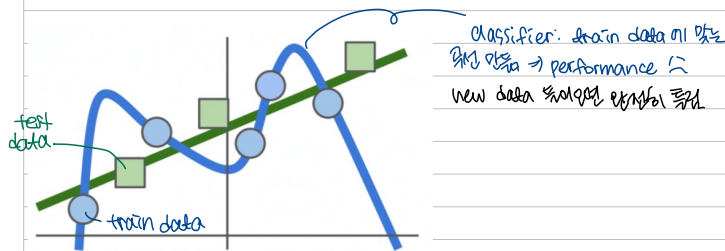
## Question about loss function

$f(x, W) = Wx$

$L = \frac{1}{N}\sum_{i=1}^{N}\sum_{j\neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$

Q. Suppose that we found a W such that L=0.
Is this W unique?
→ 유일한 W가 아니다.
ex) 2W is also has L=0.

If there exists W & 2W → 스코어 차이의 margin 또한 2배 차이 이 0

margin의 이미 >1: 괜찮은 상황. 기. 그러므로 loss는 똑같이 0이 된다.



classifier: train data에 맞는 분류 것보다 → performance ↑
new data 관점에서 일반화의 측면

test data

train data

우리는 training data가 아닌 test data에 적합한 model 찾는 경우.
→ training data의 loss에만 집중 ✗
test data 에서의 performance 집중 !! ∴ Regularization

---

Loss function now has:
- Data Loss
- Regularization Loss
+ Lambda (Hyperparameter)

$$L(W) = \frac{1}{N}\sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \boxed{\lambda R(W)}$$

추가 항목이 첨부가!

① Data loss term:
model predictions Should match
training data

② Regularization:
model should be "simple",
so it works on test data.

"simple": 해석 문제, 생성이 따라 줌이 - Occam's Razor: "Among competing hypothesis, the simple is the best"

일반적으로 더 단순한 것을 선택해야 함

→ 단순할 수록 미래 해석 설명가능성 ↑

∴ Regularization penalty

Regularization → R [ 1) 모델 의 복잡해지는것 방지
                    2) 극단 penalty 부여

$\lambda$ = regularization strength (hyperparameter)

## Regularization

Weight Won 대한 Euclidean Norm
(squared Norm) (or, ½ squared Norm)
→ won 대한 penalty 부여하기.
(L1 Regularization은
작은 won에 penalty 부여)

$$L = \frac{1}{N}\sum_{i=1}^{N}\sum_{j\neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

**L2 regularization**   $R(W) = \sum_k \sum_l W_{k,l}^2$
L1 regularization   $R(W) = \sum_k \sum_l |W_{k,l}|$
Elastic net (L1 + L2)   $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$
Max norm regularization (might see later)
Dropout (will see later)
Fancier: Batch normalization, stochastic depth ✓

→ 앞으로의 차시에서 ... 배워 나감.

∴ Regulation: training set의 overfitting 방지하기 위해
penalty 부여하는 방법!

36/85

## L2 Regularization (Weight Decay)

$$x = [1,1,1,1] \qquad R(W) = \sum_k \sum_l W_{k,l}^2$$

$$\begin{cases} w_1 = [1,0,0,0] \\ w_2 = [0.25, 0.25, 0.25, 0.25] \end{cases}$$

(If you are a Bayesian: L2 regularization also corresponds MAP inference using a Gaussian prior on W)

$$w_1^T x = w_2^T x = 1$$

*Linear classifier에서 쓰인다면, w1과 w2는 같은 결과지만 L2 Regression은 normal 더 작은 w2 선호*

Linear Classification에서 w: "얼마나 X가 output class와 관계 있는지"

(L2) 모든 X의 인자가 출력을 영향을 주길 바랄 때 사용. (W2 선호)

(L1) → 반대로. "복잡도"를 다르게 정의함 (W1 선호)
   가중치 w에, 0이 많을수록 덜 복잡함. (0↑: 복잡함)

Spave solution 선호: W의 인자 중 대부분이 0가 될것임.

---

SVM: 스코어 '자체'에 대해 관심 X. 상대적인 채널갔다 관심있음.
← Softmax: 스코어 값자체'에 대해 관심. What those scores means?
   → 정규화를 통해 확률분포로 만들겠다 0

### Softmax Classifier (Multinomial Logistic Regression)

*Probability Distribution*

scores = unnormalized log **probabilities** of the classes.

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \qquad \text{where} \qquad s = f(x_i; W)$$

*Probabilities 0~1 → sum=1*

Softmax Function (positive value)

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

$$L_i = -\log P(Y = y_i | X = x_i)$$

| cat | **3.2** |
| car | 5.1 |
| frog | -1.7 |

in summary:  $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

∴ Probability가 1에 가까울수록 정답에 가까운 값. (good)
   (→ loss: 정답일수록 작은 값. (bad))
   → 음수를 취해야 함. (minus log of true probabilities.)

### Softmax Classifier (Multinomial Logistic Regression)

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

unnormalized probabilities

| cat | **3.2** | | **24.5** | | **0.13** |
| car | 5.1 | | 164.0 | | 0.87 |
| frog | -1.7 | | 0.18 | | 0.00 |

loss function exp  →  to sum=1 normalize

*정렬 20이 됨으로 만들기*

L_i = -log(0.13) = 0.89

unnormalized log probabilities          probabilities

**Q.** what is the min/max possible loss $L_i$?

-min: 0  (log(1)=0)  answer class.        { answer class: 1
-max: ∞  (log(0+) = ∞)                     {   ~        : 0

loss 구할때 minus (log) 취해서 양수된다!

**Q.** Usually at initialization w is small so as $s \approx 0$. what is the loss?
$-\log(1/c) = \log C$
→ log C가 아니면 뭔가 잘못됨..

"How Badness the W is?"
즉 일마나 나쁜지 를 다룬.



**matrix multiply + bias offset**

| 0.01 | -0.05 | 0.1 | 0.05 |
|------|-------|-----|------|
| 0.7 | 0.2 | 0.05 | 0.16 |
| 0.0 | -0.45 | -0.2 | 0.03 |

$W$

| -15 |
|-----|
| 22 |
| -44 |
| 56 |

+

| 0.0 |
|-----|
| 0.2 |
| -0.3 |

$b$

$x_i$

$y_i$ [2]

**hinge loss (SVM)**

| -2.85 |
|-------|
| 0.86 |
| 0.28 |

max(0, -2.85 - 0.28 + 1) +
max(0, 0.86 - 0.28 + 1)
=
**1.58**

**cross-entropy loss (Softmax)**

| -2.85 | 0.058 | 0.016 |
|-------|-------|-------|
| 0.86 | 2.36 | 0.631 |
| 0.28 | 1.32 | 0.353 |

exp → normalize (to sum to one) → - log(0.353) = **0.452**

---

Softmax      vs      [SVM]

very high score ⟶ always consider datapoint (to get better)
low score

consider only thing : "margin"
datapoint → just only we to get value.

---

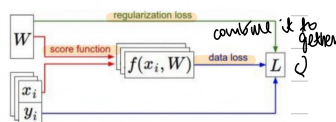## Recap

### How do we find the best W?
→ Optimization !

- We have some dataset of (x,y)
- We have a **score function:**   $s = f(x; W) \stackrel{e.g.}{=} Wx$
- We have a **loss function**: How quantity bad the W is.

$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$   **Softmax**

$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$   **SVM**

$L = \frac{1}{N}\sum_{i=1}^N L_i + R(W)$   **Full loss**



combine it to gether

---

Setting W → "iterative." improve it every time

**Strategy #1. Random Search**
Really Bad Algorithm !!

**Strategy #2. Follow the slope**
where is the slop of this? → go ahead little bit..
Get well! General strategy.

1-dim func :   $\frac{df(x)}{dx} = \lim_{h\to 0} \frac{f(x+h) - f(x)}{h}$

multi-dim func:  [gradient] is the vector of (partial derivatives) along each dimension.

tell us what is the function f that we move in backward direction.

---

| current W: | W + h (first dim): | gradient dW: |
|-----------|---------------------|--------------|
| [0.34, | [0.34 + **0.0001**, | [-2.5, |
| -1.11, | -1.11, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | (1.25322 - 1.25347)/0.0001 |
| 0.55, | 0.55, | = -2.5 |
| 2.81, | 2.81, | $\frac{df(x)}{dx} = \lim_{h\to 0}\frac{f(x+h)-f(x)}{h}$ |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?,...] |
| 0.33,...] | 0.33,...] | |
| **loss 1.25347** | **loss 1.25322** | |

Super slow !!

use math

This is silly. The loss is just a function of W:

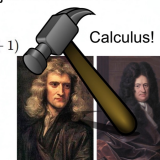$L = \frac{1}{N}\sum_{i=1}^N L_i + \sum_k W_k^2$
$L_i = \sum_{j\neq y_i} \max(0, s_j - s_{y_i} + 1)$
$s = f(x; W) = Wx$

want $\nabla_W L$

Use calculus to compute an
**analytic gradient**

Calculus!

| current W: | | gradient dW: |
|-----------|---|--------------|
| [0.34, | dW = ... | [-2.5, |
| -1.11, | (some function | 0.6, |
| 0.78, | data and W) | 0, |
| 0.12, | | 0.2, |
| 0.55, | | 0.7, |
| 2.81, | | -0.5, |
| -3.1, | | 1.1, |
| -1.5, | | 1.3, |
| 0.33,...] | | -2.1,...] |
| **loss 1.25347** | | |

go directly (one step)

way more efficient. fast

(In practice: Always we analytic gradient,
but check implementation with numerical gradient.)
⇒ "gradient check"

very useful debugging strategy

super slow..
but super useful for debugging...

code

```
# Vanilla Gradient Descent

while True:
  weights_grad = evaluate_gradient(loss_fun, data, weights)
  weights += - step_size * weights_grad # perform parameter update
```

hyper parameter

(learning rate, step size : 처음에 설정하는 parameter.)

## Stochastic Gradient Descent (SGD)

many data → so slow
rather than computing
entire set,
we minibatch
small → estimate of
the small sum.

update
gradient
with this
estimate!

code →

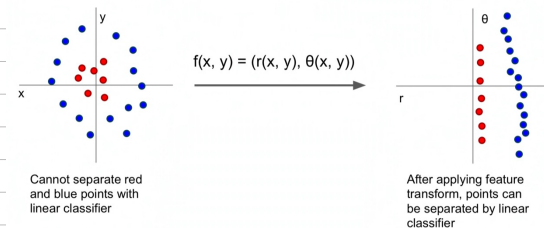$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive
when N is large!

Approximate sum
using a **minibatch** of
examples
32 / 64 / 128 common

```
# Vanilla Minibatch Gradient Descent

while True:
  data_batch = sample_training_data(data, 256) # sample 256 examples
  weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
  weights += - step_size * weights_grad # perform parameter update
```

## Image Features

→ Linear classifier 보다 앞 ∴ Multi-Modality
① calculate feature
→ vector → to linear classifier input.

Motivation → Feature transform

## Image Features: Motivation



f(x, y) = (r(x, y), θ(x, y))

Cannot separate red
and blue points with
linear classifier

After applying feature
transform, points can
be separated by linear
classifier

NN이 뜨기 전, HOG (Local orientation edges 특징) 이용