

# オンライン学習

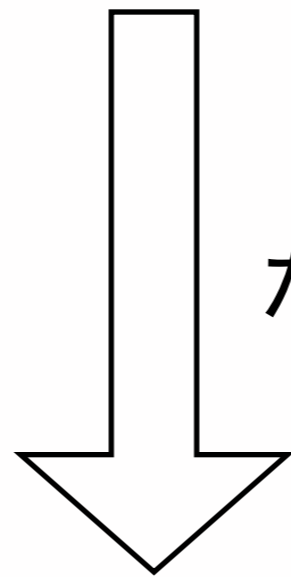
渡辺太郎

taro.watanabe at nict.go.jp



<https://sites.google.com/site/alaginmt2014/>

機械翻訳について勉強したい。



なるべく良い翻訳..

I want to study about machine translation.  
I need to master machine translation.  
machine translation want to study.

infobox )

infobox buddhist

道元は鎌倉時代初期のである。

道元（どうげん）は、鎌倉時代初期の禅僧。

曹洞禅の開祖

曹洞宗の開祖。

その生涯には kigen 名がある。

晩年に希玄という異称も用いた。

一般には宗と呼ばれることによって尊称は高僧がある。

同宗旨では高祖と尊称される。

死後にといった仏所伝灯国師、 joyo-daishi である。

諡は、仏性伝東国師、承陽大師\_(僧)。

一般には道元禅師と呼ばれる。

一般には道元禅師と呼ばれる。

また～686の普及についての修行を tooth brushing、大峰、食事作法

cleaning として、日本にしている。

日本に歯磨き洗面、食事の際の作法や掃除の習慣を広めたといわれ

れる。

# エラー

- 探索エラー: スコアの高い翻訳を出すのに失敗
- モデルエラー: スコアの高い翻訳が誤っている
- 学習データの問題: 小さい、異なる
- 手軽な対処: 最適化(チューニング)

# どうしまししょう?

$f$  = 機械翻訳について勉強したい

$$\log Pr(\phi|e) \log Pr(e) \log Pr(f, \alpha|\phi)$$

I want to study about machine translation

I need to master machine translation

machine translation want to study

I don't want to learn anything

-2	-3	-4	-9
-3	-4	-4	-11
-2	-5	-1	-8
-5	-2	-3	-10

kbestを正しく並び替

えるように重みを学習

0.5×-2	0.4×-3	0.2×-4	-3.0
0.5×-3	0.4×-4	0.2×-4	-3.9
0.5×-2	0.4×-5	0.2×-1	-3.2
0.5×-5	0.4×-2	0.2×-3	-3.9

# 重み付け

$$\hat{e} = \arg \max_e Pr(\mathbf{f}, \alpha | \phi, \mathbf{e})^{0.2} Pr(\phi | \mathbf{e})^{0.5} Pr(\mathbf{e})^{0.4}$$

- より一般化:

$$\begin{aligned}\hat{e} &= \arg \max_e \frac{\sum_d \exp(\mathbf{w}^\top \mathbf{h}(\mathbf{f}, d, \mathbf{e}))}{\sum_{e', d'} \exp(\mathbf{w}^\top \mathbf{h}(\mathbf{f}, d', e'))} \\ &\approx \arg \max_{\langle \mathbf{e}, d \rangle} \mathbf{w}^\top \mathbf{h}(\mathbf{f}, d, \mathbf{e})\end{aligned}$$

- 複数の素性 $\mathbf{h}(\mathbf{e}, d, \mathbf{f})$ をlog-linearに組み合わせ  
最適化 = 最適な $\mathbf{w}$ を学習



# MTパイプライン

kyoto-train.{ja,en}

対訳データ

大量(低品質?)

翻訳モデル

言語モデル

デコーダ

最適化の問題

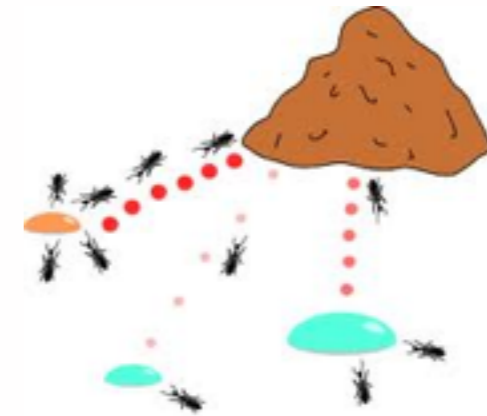
重みパラメータ

(少量?)高品質

対訳データ

kyoto-tune.{ja,en}

# 最適化



$$\begin{aligned}\hat{\boldsymbol{w}} &= \arg \min_{\boldsymbol{w} \in \mathcal{W}} \mathbb{E}_{Pr(F, E)} [\ell(F, E; \boldsymbol{w})] \\ &= \arg \min_{\boldsymbol{w} \in \mathcal{W}} \ell(F, E; \boldsymbol{w}) + \lambda \Omega(\boldsymbol{w})\end{aligned}$$

- ある損失関数  $\ell$  を仮定し、対訳データ  $(F, E)$  に対するリスクを最小化□□
- 真の分布は未知なため、正則化 ( $\Omega$ ) された経験リスクを最小化



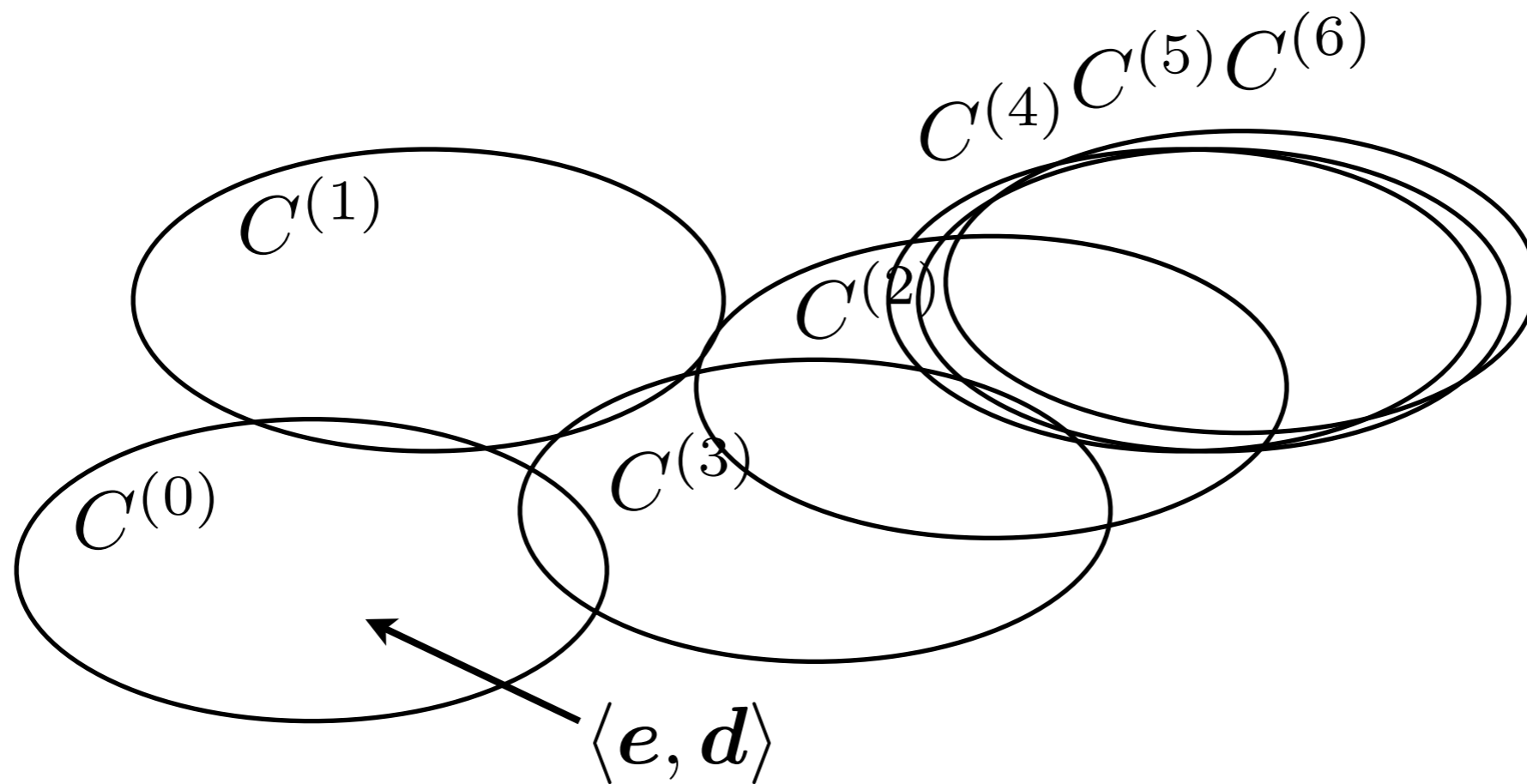
# kbest近似最適化

```
1: procedure LEARN( $\langle F, E \rangle$ )
2:    $C = \{\}$ 
3:   for  $t = 1 \dots T$  do
4:      $w$ でデコード、kbestリストを生成
5:     kbestリストを $C$ へ結合
6:      $w^{(t+1)} = \arg \min_w \ell(F, E, C; w) + \lambda \Omega(w)$ 
7:   end for
8: end procedure
```

(Och and Ney, 2002)

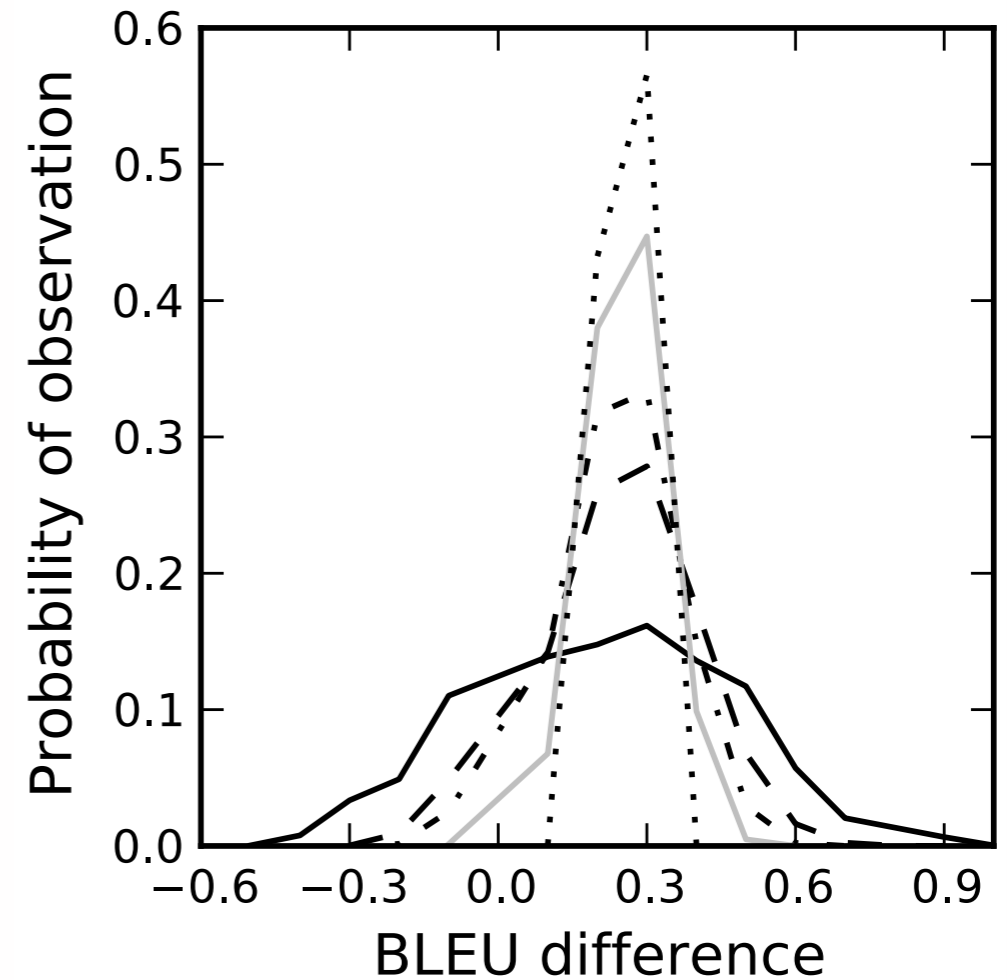
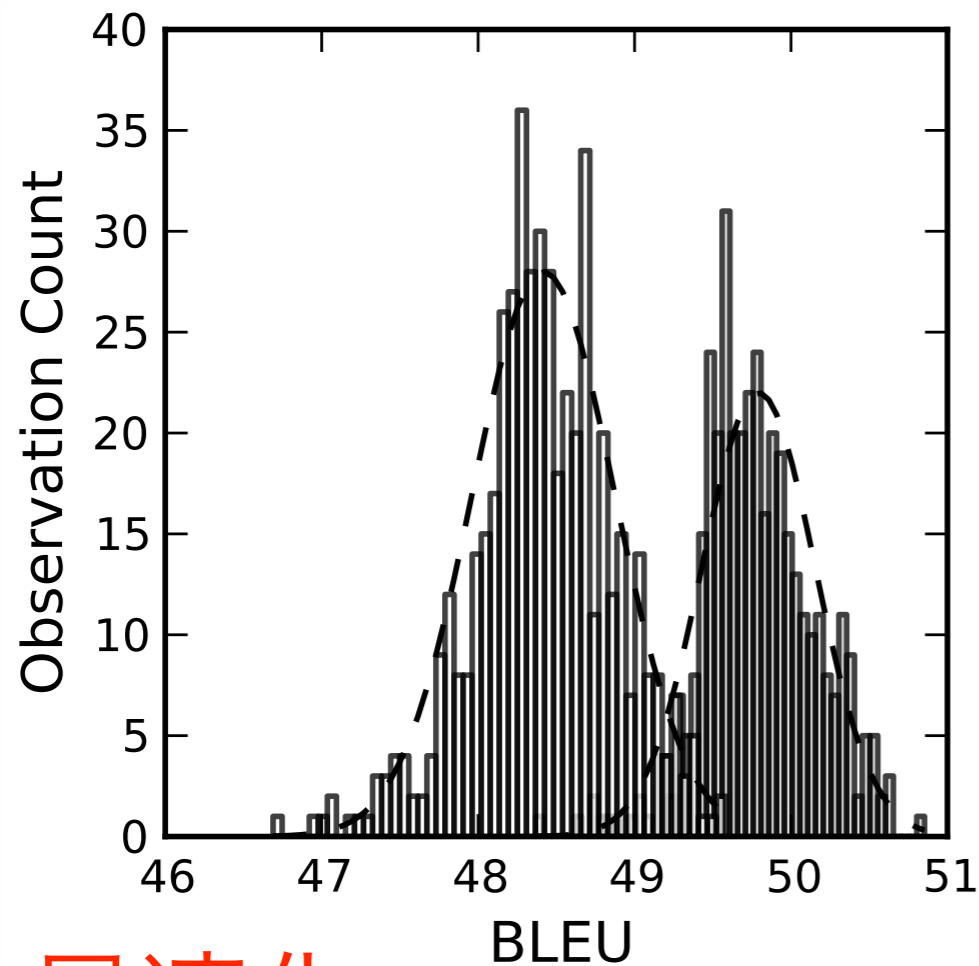
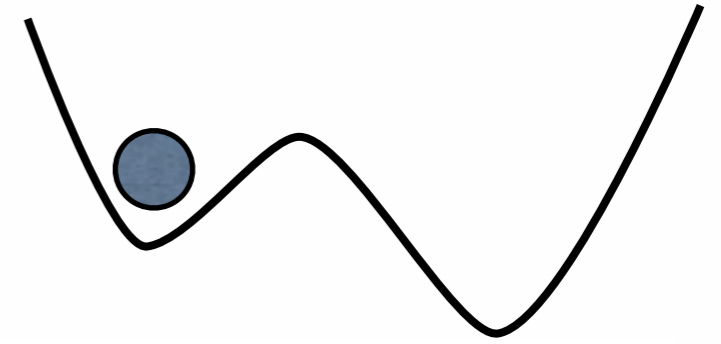
- “k-best結合”以外の近似?

# kbest近似最適化



- 各繰り返しでk-bestを結合しつつ学習
- 収束はするけど... 局所解

# 局所解



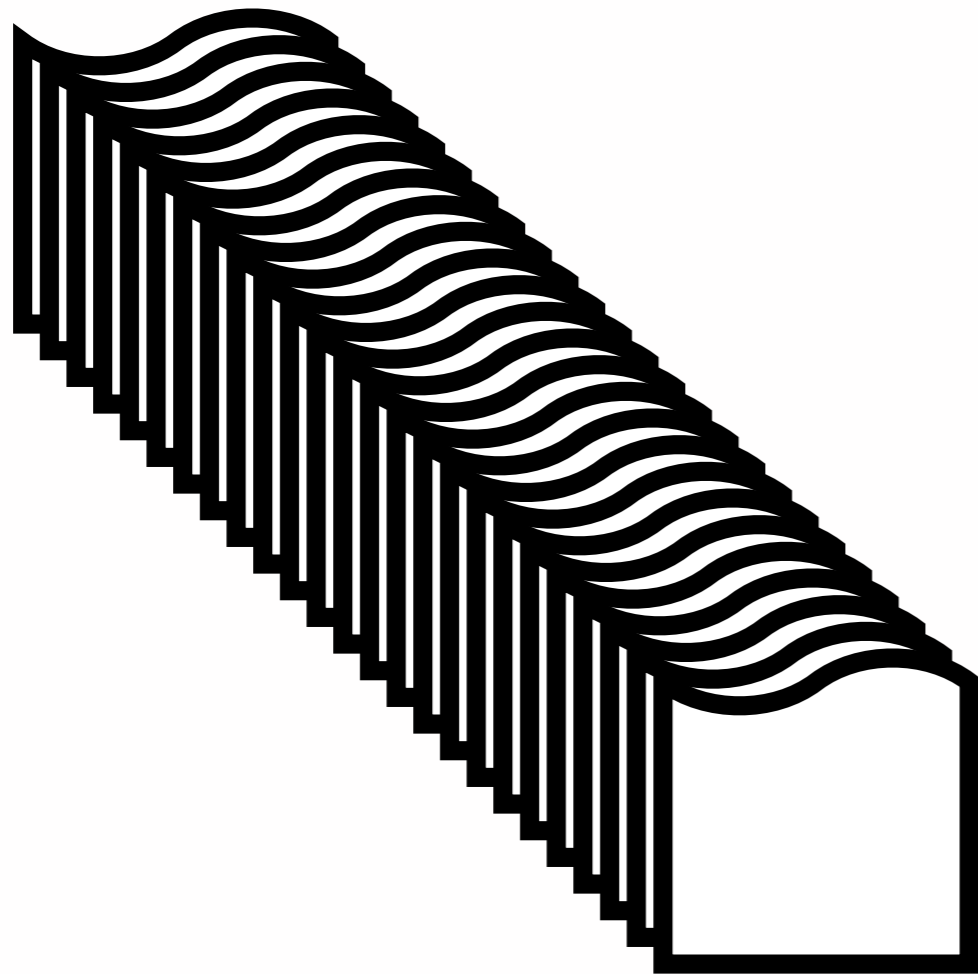
(Clark et al., 2011)

## 最適化

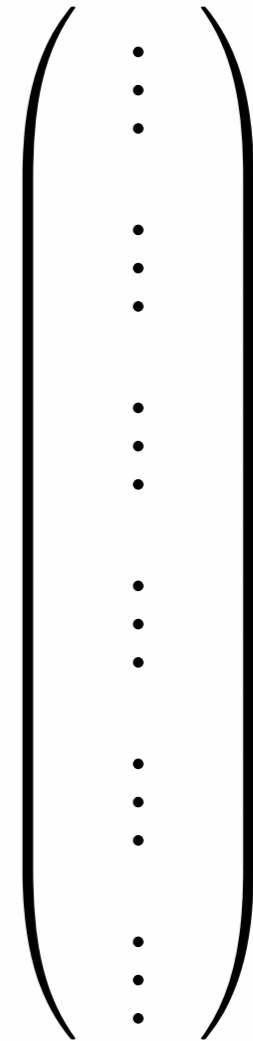
- ~~MERT~~ は、パラメータの初期値に大きく依存
- ランダムな初期値から複数 ~~MERT~~

## 最適化

# 大規模化の問題

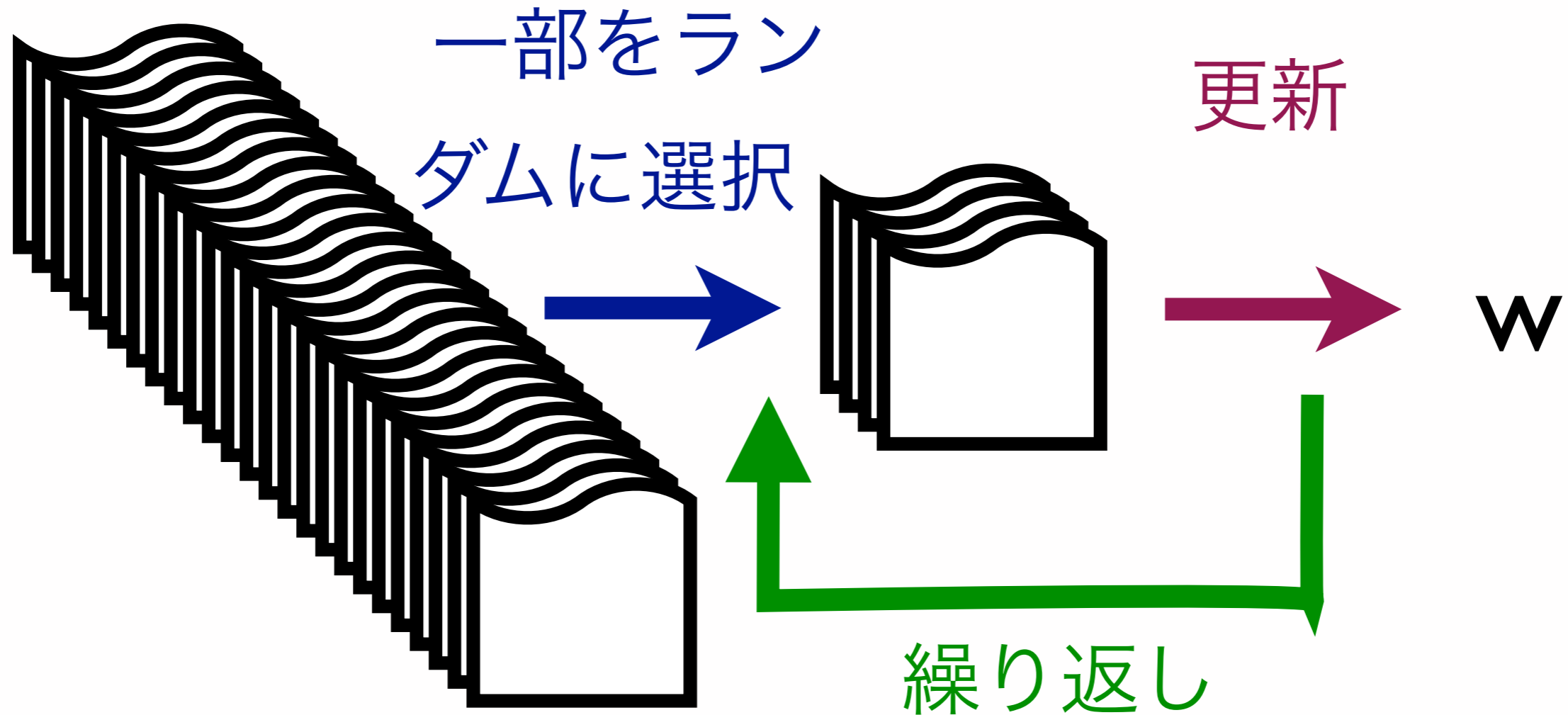


$$h(f, d, e) =$$



- データが多い、高次元な素性ベクトル
- 効率のよい学習手法?

# どうせ近似するならば...



- 学習データを近似
- 非常に簡単なアルゴリズムで実現

# オンライン学習

```
1: procedure ONLINELEARN( $\langle F, E \rangle = \left\{ \langle f^{(i)}, e^{(i)} \rangle \right\}_{i=1}^N$ )
2:    $w^{(1)} \leftarrow \emptyset$ 
3:   for  $t \in \{1, \dots, T\}$  do
4:      $\langle \tilde{F}^{(t)}, \tilde{E}^{(t)} \rangle \subseteq \langle F, E \rangle$ 
5:      $\tilde{C}^{(t)} \leftarrow \text{GEN}(\tilde{F}^{(t)}, w^{(t)})$ 
6:      $w^{(t+1)} \leftarrow \arg \min_{w \in \mathcal{W}} \ell(\tilde{F}^{(t)}, \tilde{E}^{(t)}, \tilde{C}^{(t)}; w) + \lambda \Omega(w)$ 
7:   end for
8:   return  $w^{(T+1)}$ 
9: end procedure
```

ランダムにサンプル  
リング: mini batch

デコード

更新

- 自然言語処理ではよく使われる: 形態素解析、構文解析 etc.
- 当然、機械翻訳にも使えます (Watanabe et al., 2007)

# 問題: BLEU



$e_1$  I want to study about machine translation

$e_2$  I need to master machine translation

$e_3$  machine translation want to study

$e_4$  I don't want to learn anything

- k-bestのランキングに基づく最適化
- コーパス単位のBLEU  $\neq$  文単位のBLEUの平均
- 文単位では、正しく学習するのは困難

# BLEUの近似

$$\text{GEN}(f^{(s)}, w)$$

$$e^{(1)}, \dots, \begin{pmatrix} c_1^{(s)} \\ \vdots \\ c_i^{(s)} \\ \vdots \\ c_K^{(s)} \end{pmatrix}, \dots, e^{(S)}$$

- 今までの各文に対するBLEUの統計量を保存(I-bestあるいはoracle) (Watanabe et al., 2007)



# 減衰によるBLEUの近似

$$\mathbf{b} \leftarrow 0.9 \times (\mathbf{b} + \mathbf{c}(e))$$

$$l \leftarrow 0.9 \times (l + |\mathbf{f}|)$$

$$B(e) = (l + |\mathbf{f}|) \times \text{Bleu}(\mathbf{b} + \mathbf{c}(e))$$

$$\hat{e}^{(s)} = \underset{e}{\operatorname{argmax}} -B(e) + \mathbf{w}^\top \mathbf{h}(\mathbf{f}^{(s)}, e)$$

$$\dot{e}^{(s)} = \underset{e}{\operatorname{argmax}} +B(e) + \mathbf{w}^\top \mathbf{h}(\mathbf{f}^{(s)}, e)$$

- sentence-BLEUに対して、今までのBLEUの履歴 ( $\times 0.9$ )を加える (Chiang et al., 2008)

# パーセプトロン

```
1: procedure PERCEPTRON( $\langle F, E \rangle = \left\{ \langle f^{(i)}, e^{(i)} \rangle \right\}_{i=1}^N$ )
2:    $w^{(1)} \leftarrow 0$ 
3:   for  $t \in \{1 \dots T\}$  do
4:      $\langle f, e \rangle \sim \langle F, E \rangle$ 
5:      $\tilde{c} \leftarrow \text{GEN}(f, w^{(t)})$ 
6:      $\langle \hat{e}, \hat{d} \rangle \in \tilde{c}$ 
7:      $\langle e^*, d^* \rangle \in \tilde{o} \subseteq \tilde{c}$ 
8:     if  $\langle e^*, d^* \rangle \neq \langle \hat{e}, \hat{d} \rangle$  then
9:        $w^{(t+1)} \leftarrow w^{(t)} + h(f, e^*, d^*) - h(f, \hat{e}, \hat{d})$ 
10:    end if
11:  end for
12:  return  $w^{(T+1)}$ 
13: end procedure
```

一文だけサンプル

I-best

オラクル翻訳

間違ったらペナルティ

最後のパラメータ or 平均値

# パーセプトロン

$e_3$  machine translation want to study

$e_2$  I need to master machine translation

$e_1$  I want to study about machine translation

$e_4$  I don't want to learn anything

$$\arg \min_w \max \{0, 0 - (w^\top h(f, e_1) - w^\top h(f, e_3))\}$$

- $e_1$ のスコアが、 $e_3$ のスコアよりも、大きくなるように、 $w$ を更新

# hinge損失

$e_3$  machine translation want to study

$e_2$  I need to master machine translation

$e_1$  I want to study about machine translation

$e_4$  I don't want to learn anything

$$\arg \min_w \max \{0, 1 - (w^\top h(f, e_1) - w^\top h(f, e_3))\}$$

- $e_1$ のスコアが、 $e_3$ のスコアよりも、**1以上**大きくなるように、 $w$ を更新

# MIRA

$$\arg \min_w \lambda \frac{1}{2} \|w - w^{(t)}\|^2 + \Delta \text{error}(e_1, e_3) - (w^\top h(f, e_1) - w^\top h(f, e_3))$$

(Crammer et al., 2006)

- 以前のパラメータ  $w^{(t)}$  との差分を小さくしつつ、 $e_1$  のスコアが、 $e_3$  のスコアよりも、**翻訳のエラー以上**大きくなるように、 $w$  を更新

# MIRA

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \alpha^{(t)} \Delta \mathbf{h}(\mathbf{f}, \mathbf{e}_1, \mathbf{e}_3)$$

$$\alpha^{(t)} = \min \left\{ \frac{1}{\bar{\lambda}}, \frac{\Delta \text{error}(\mathbf{e}_1, \mathbf{e}_3) - \mathbf{w}^{(t)\top} \Delta \mathbf{h}(\mathbf{f}, \mathbf{e}_1, \mathbf{e}_3)}{\|\Delta \mathbf{h}(\mathbf{f}, \mathbf{e}_1, \mathbf{e}_3)\|^2} \right\}$$

$$\Delta \mathbf{h}(\mathbf{f}, \mathbf{e}_1, \mathbf{e}_3) = \mathbf{h}(\mathbf{f}, \mathbf{e}_1) - \mathbf{h}(\mathbf{f}, \mathbf{e}_3)$$

- 非常に簡単な更新:  $\alpha^{(t)}$ により更新の量を調整
- $\alpha^{(t)}=1$ : パーセプトロン

他にも、CWやAROWなど様々な更新式

# 大きいマージン

$e_3$  machine translation want to study

$e_2$  I need to master machine translation

$e_1$  I want to study about machine translation

$e_4$  I don't want to learn anything

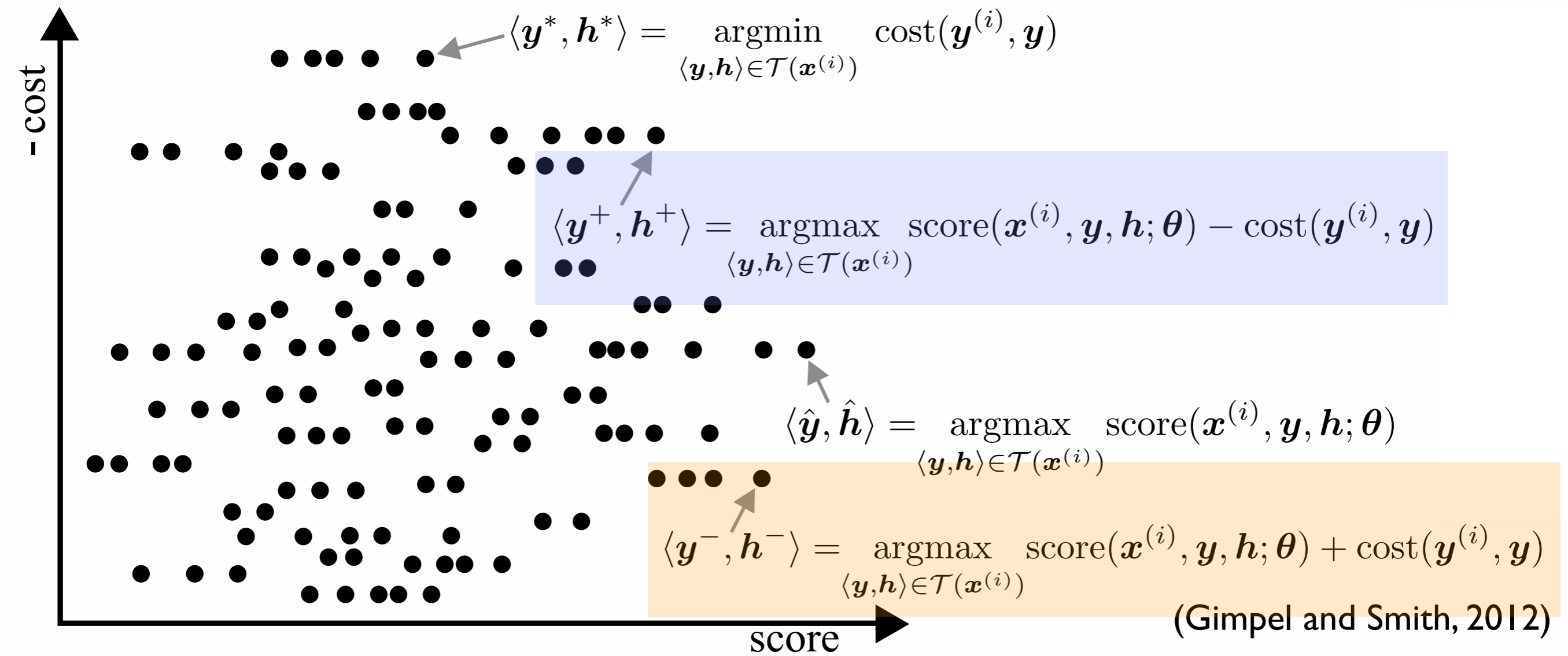
$$\Delta \text{error}(e_1, e_3) - w^{(t)\top} \Delta h(f, e_1, e_3)$$

あるいは

$$\Delta \text{error}(e_2, e_3) - w^{(t)\top} \Delta h(f, e_2, e_3)$$

- 正しい翻訳と誤った翻訳のうち、誤りの差が大きく、かつ、スコアの差が大きいペアを選択

(Chiang et al., 2008; Chiang et al., 2009)





# SGD

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \eta^{(t+1)} \left( \Delta \ell(\tilde{F}^{(t)}, \tilde{E}^{(t)}, \tilde{C}^{(t)}; \boldsymbol{w}) + \lambda \Delta \Omega(\boldsymbol{w}) \right)$$

- 勾配が計算できるなら、hinge損失以外でも使える
- 学習率 $\eta$ は、学習が進むにつれて、減衰

# AdaGrad

$$\mathbf{g}^{(t+1)} \leftarrow \mathbf{g}^{(t)} + \Delta\ell(\tilde{F}^{(t)}, \tilde{E}^{(t)}, \tilde{C}^{(t)}; \mathbf{w})^2$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \frac{\eta_0}{\sqrt{\mathbf{g}^{(t+1)}}} \left( \Delta\ell(\tilde{F}^{(t)}, \tilde{E}^{(t)}, \tilde{C}^{(t)}; \mathbf{w}) + \lambda\Delta\Omega(\mathbf{w}) \right)$$

(Duchi et al., 2011)

- 各素性毎に、学習率を自動的に調整
- 他にも、RDA、AdaDec、AdaDeltaなど

# 並列化

```
1: procedure PARALLELEARN( $\langle F, E \rangle$ )
2:    $\{\langle F_1, E_1 \rangle, \dots, \langle F_S, E_S \rangle\} \leftarrow \text{SPLIT}(\langle F, E \rangle)$ 
3:   for  $s \in \{1, \dots, S\}$  parallel do
4:      $w_s \leftarrow \text{ONLINELEARN}(\langle F_s, E_s \rangle)$ 
       あるいは  $\text{BATCHLEARN}(\langle F_s, E_s \rangle)$ 
5:   end for
6:    $w \leftarrow \text{mix}(\{w_1, \dots, w_S\})$ 
7:   return  $w$ 
8: end procedure
```

shardへ分割

各shardで学習

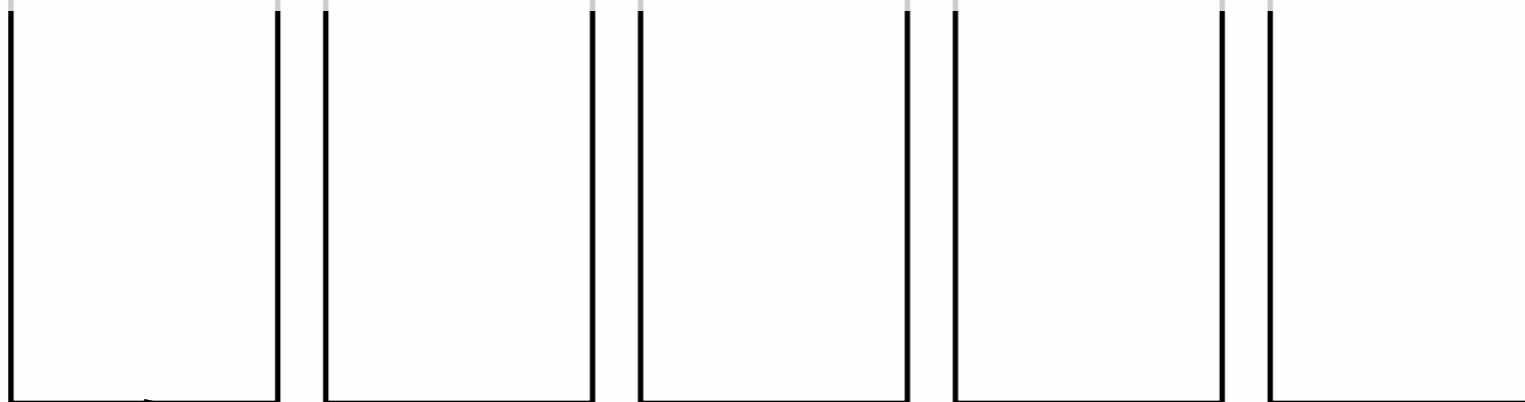
混合

- 各shard独立に学習されるため、全体で最適解が保証されない

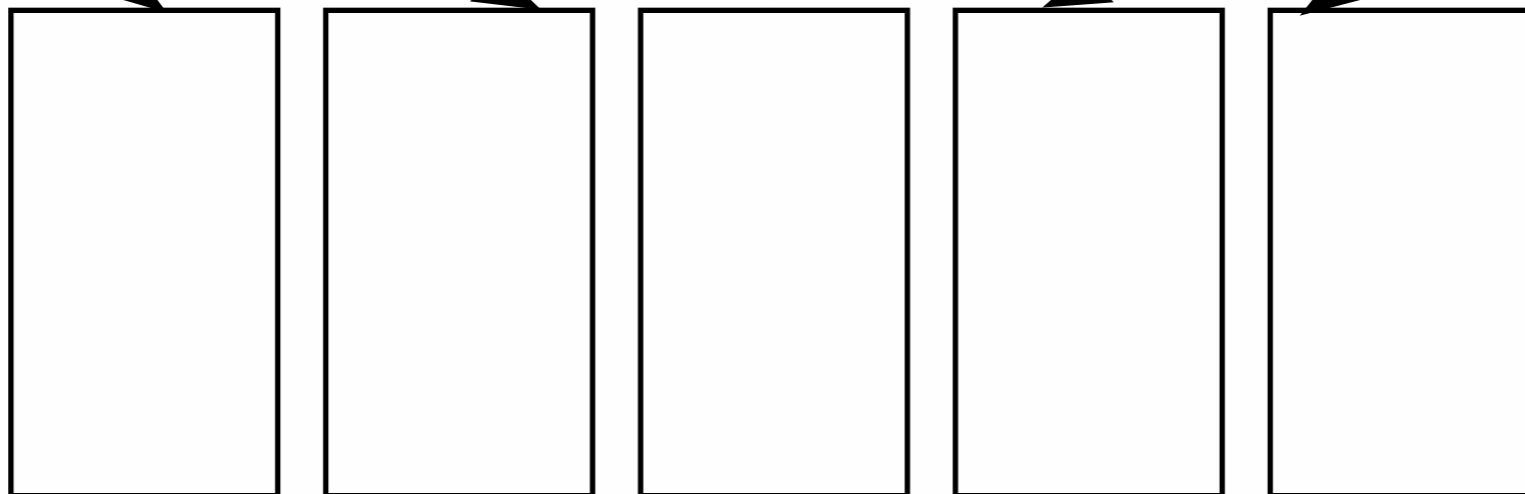
# 同期更新

shard1 shard2 shard3 shard4 shard5

t



繰り返し  
毎に混合

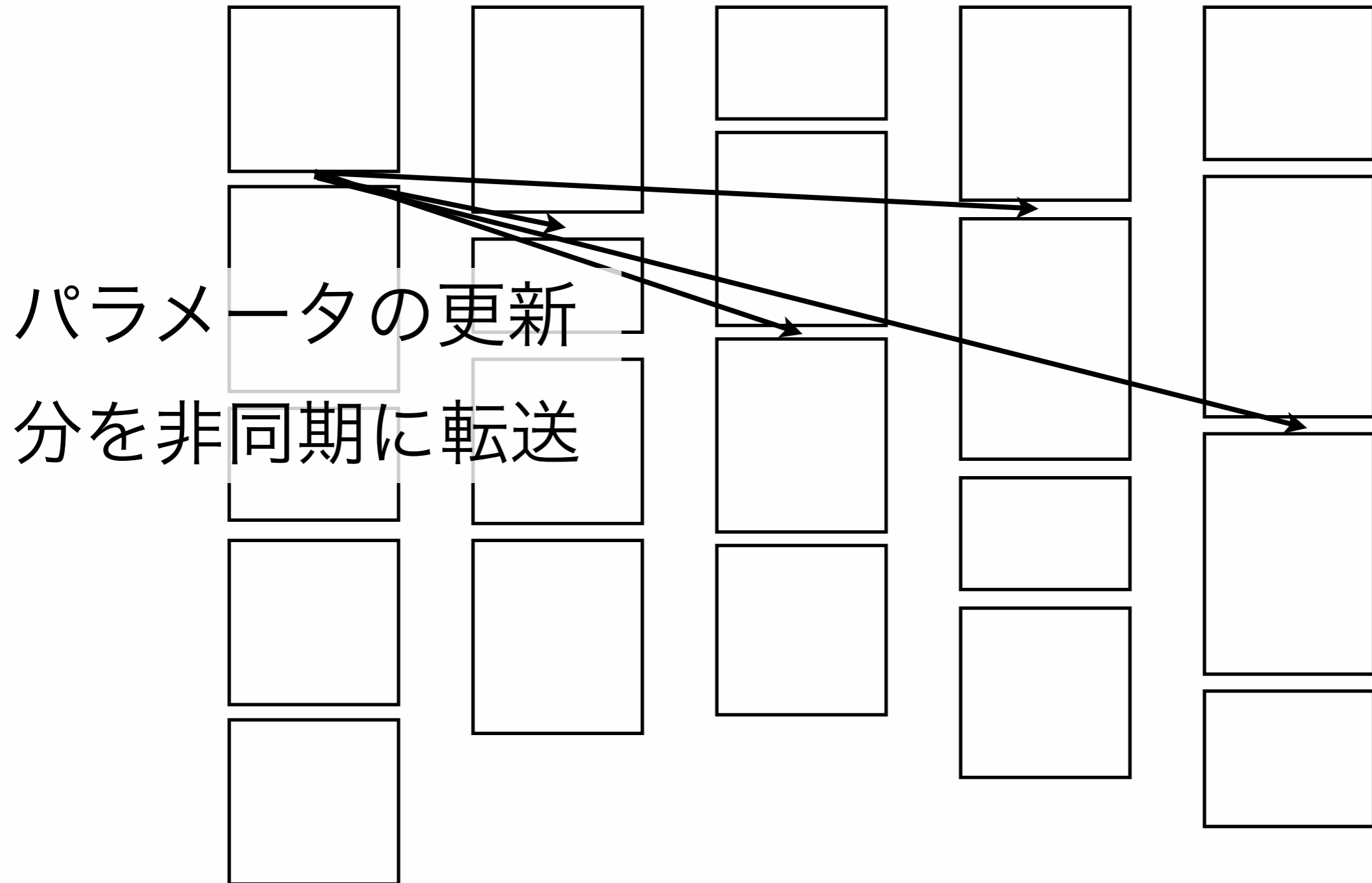


t + 1



(McDonald et al., 2010)

# 非同期更新



# まとめ

- オンライン学習: 大規模データ、高次元索性
- パーセプトロン、MIRA、SGD、AdaGrad
- 並列化による効率化

# 参考文献

- David Chiang, Yuval Marton, and Philip Resnik. 2008. Online large-margin training of syntactic and structural translation features. In *Proc. of EMNLP 2008*, pages 224-233, Honolulu Hawaii, October.
- David Chiang, Kevin Knight, and Wei Wang. 2009. |1,001 new features for statistical machine translation. In *Proc. of NAACL-HLT 2009*, pages 218-226, Boulder, Colorado, June
- Jonathan H. Clark, Chris Dyer, Alon Lavie, and Noah A Smith. 2011. Better hypothesis testing for statistical machine translation: Controlling for optimizer instability. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 176-181, Portland, Oregon, USA, June.
- Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. 2006. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7:551-585 March.
- John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121-2159, July.

# 参考文献

- Kevin Gimpel and Noah A. Smith. 2012. Structured ramp loss minimization for machine translation. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 221-231, Montréal, Canada, June.
- Ryan McDonald, Keith Hall, and Gideon Mann. 2010. Distributed training strategies for the structured perceptron. In *Proc. of NAACL-HLT 2010*, pages 456-464, Los Angeles California, June.
- Franz Josef Och and Hermann Ney. 2002. Discriminative training and maximum entropy models for statistical machine translation. In *Proc. of ACL 2002*, pages 295-302 Philadelphia, Pennsylvania, USA, July.
- Taro Watanabe, Jun Suzuki, Hajime Tsukada, and Hideki Isozaki. 2007. Online Large-Margin Training for Statistical Machine Translation. In *Proc. of EMNLP-CoNLL 2007*, pages 764-773, Prague, Czech Republic, June.