

today: FW & BCFW for structured SVM!

recall: batch FW on

minus dual objective  $-d(\alpha)$

$$\min_{\alpha \in A_{\text{reg}}} \frac{\lambda \|A\alpha\|^2 - b^T \alpha}{2}$$

$$(\nabla d(\alpha))_{i,y} = -\frac{1}{n} H(y_i; w(\alpha)) \quad w(\alpha) = A\alpha \\ = (A^T A - b)_{i,y}$$

FW step on  $M = \sum_{i=1}^n \Delta_{y_i \alpha_i}$

loss-augmented inference for each  $i$   $\tilde{y}_i^{(t)} = S_{y_i^{(t)}} \text{ where } \tilde{y}_i^{(t)} = \arg \max_{\tilde{y} \in \mathcal{Y}_i} H_i(\tilde{y}; w^{(t)})$

$$\alpha^{(t)} \xrightarrow{A} w^{(t)}$$

$$\text{from } \alpha^{(t+1)} = (1-\gamma) \alpha^{(t)} + \gamma S^{(t)}$$

$$\text{recall primal } p(w) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n H_i(w)$$

$$p(w, \tilde{y}(w)) \quad \frac{\partial p}{\partial w^{(t)}} = \lambda w - \frac{1}{n} \sum_{i=1}^n \nabla_i(\tilde{y}_i^{(t)})$$

↑  
a subgradient of  $p$  at  $w^{(t)}$

$$w^{(t+1)} = (1-\gamma) w^{(t)} + \gamma \sum_{i=1}^n \nabla_i(\tilde{y}_i^{(t)})$$

subgradient step with step-size  $\beta$

$$w^{(t+1)} = (1-\beta\lambda) w^{(t)} + \beta \sum_{i=1}^n \nabla_i(\tilde{y}_i^{(t)})$$

thus batch FW on dual with step-size  $\gamma = \beta/\lambda$   
is equivalent to batch subgradient on primal  
with stepsize  $\beta = \frac{\gamma}{\lambda}$

note: subgradient convergence proof for  $\mu$ -strongly convex obj.

$$\text{needed } \beta \leq \frac{1}{\mu} O(\frac{1}{t})$$

$$R \triangleq \max_{i,y} \|\nabla_i(y)\|_2$$

$$\text{canonical FW step-size } \gamma_t = \frac{2}{t+2} \Rightarrow \beta_t = \frac{1}{t+2}$$

$\gamma_t$  strong convexity parameter of primal

$$p(\tilde{w}^{(t)}) - p(w^*) \leq \frac{8R^2}{t+2}$$

$\tilde{w}^{(t)} \triangleq \sum_{s=0}^t \beta_s^t w^{(s)}$  (weighted average)

for any convex comb.  
of  $p$

$$\min_{S \subseteq t} p(w^{(S)}) \leq \sum_{s \in S} \beta_s^t p(w^{(s)})$$

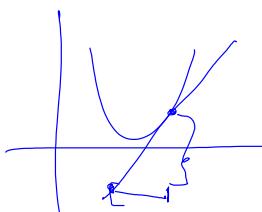
convexity of  $p$

FW analysis :  $d(\alpha^*) - d(\alpha^{(t)}) \leq \frac{2C_f}{t+2}$ ; turns out that  $C_f = \frac{4R^2}{\lambda}$  (later today)

Let's look at FW gap &  $\langle \nabla f(\alpha^{(t)}), s^{(t)} - \alpha^{(t)} \rangle = \frac{1}{n} \sum_{i=1}^n (H_i(\hat{y}_i^{(t)}; w^{(t)}) - \sum_{y \in S_i} \alpha_i(y) H_i(y; w^{(t)}) )$

$$= p(w(\alpha^{(t)})) - d(\alpha^{(t)})$$

(Subgradient gap of lecture 9?)



[note: not always the case;  
see e.g. CRF objective]

\* one can show, that  $\min_{S \subseteq t} C_S \leq 3 \cdot \frac{2C_f}{t+2}$

weighted average

it's also true that  $g^{FW}(\hat{\alpha}^{(t)}) \leq 3 \frac{2C_f}{t+2}$  when  $g(\cdot)$  is convex

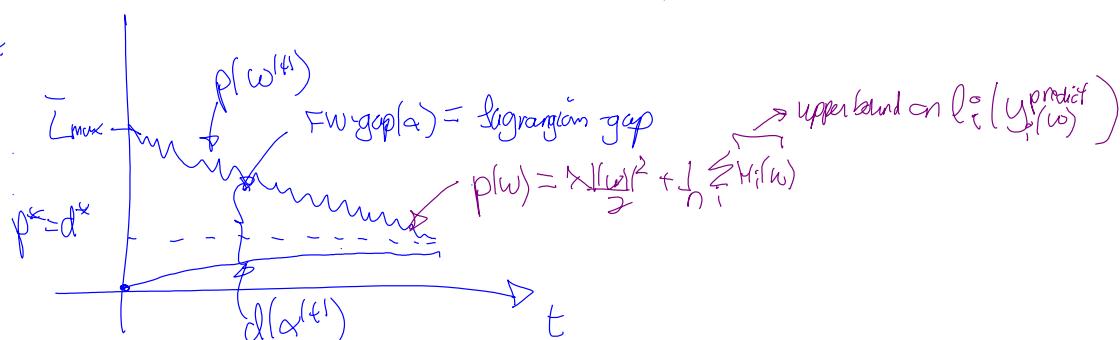
$\Lsh$  batch FW on dual for SVMstruct  $g^{FW}(\hat{\alpha}^{(t)})$

also gives  $p(\hat{w}^{(t)}) - p(w^*) = O(\frac{1}{t})$

This is case for quadratic functions

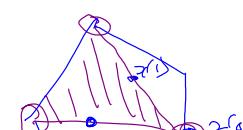
$$p(\hat{w}) = \frac{1}{n} \sum_{i=1}^n \max_y l_i(\hat{y}_i) = L_{\max}$$

$$d(\alpha^{(0)}) = 0$$



another variant of FW: FCFW "Fully-corrective FW"

algorithm: reoptimize f over convex hull  $\{S^{(u)}\}_{u=0}^t$



turns out that

(batch) FCFW is dual of SVM struct  
is equivalent to the  
constraint generation / cutting plane approach  
for the 1-stack formulation

(see Lecture 9)

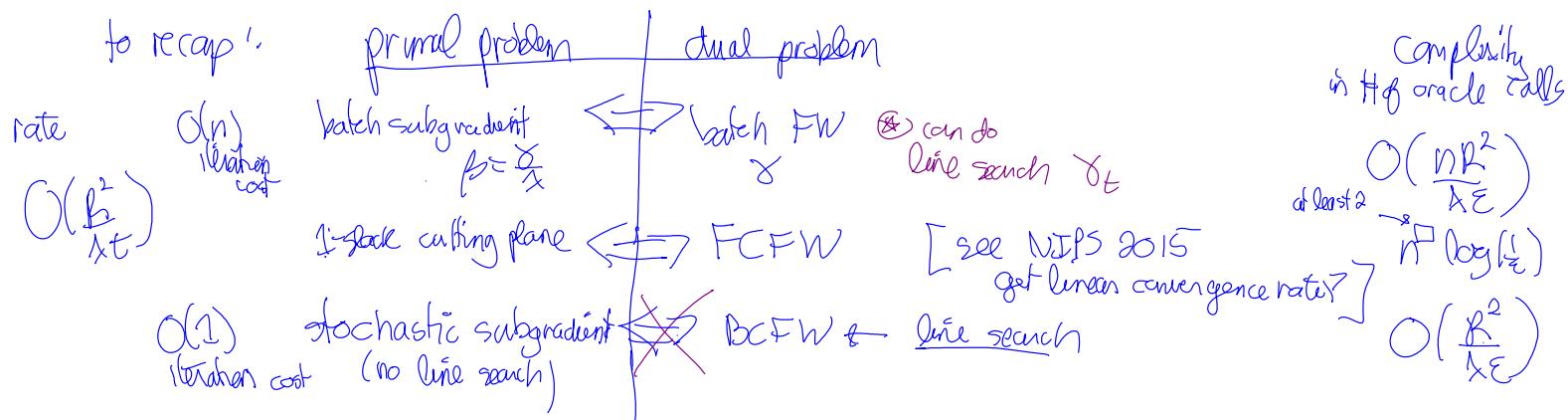
why? every  $s^{(t)}$  corresponds to  $\{y_i^{(t)}\}_{i=1}^n$

$$\alpha \in \text{convex-hull}(\{s^{(u)}\}_{u \leq t}) \Rightarrow \alpha = \sum_{u \leq t} \tilde{\alpha}_u s^{(u)}$$

$\tilde{\alpha}$  belongs to a subset of  $\Delta_{\{x_i^{(t)}\}_{i=1}^n}$

$$w = Ax = \sum_{u \leq t} \tilde{\alpha}_u A s^{(u)}$$

$\frac{1}{\lambda} \left( \sum_{i=1}^n y_i^{(t)} / \hat{y}_i^{(u)} \right)$



### pointers:

- the usual pointer for FW for structured SVM (as in previous lecture):
  - Block-Coordinate Frank-Wolfe Optimization for Structured SVMs, S. Lacoste-Julien, M. Jaggi, M. Schmidt and P. Pletscher, ICML 2013  
<http://jmlr.csail.mit.edu/proceedings/papers/v28/lacoste-julien13-supp.pdf>
- the linear convergence of FCFW is given in:
  - On the Global Linear Convergence of Frank-Wolfe Optimization Variants, S. Lacoste-Julien and M. Jaggi, NIPS 2015  
<http://arxiv.org/abs/1511.05932> | [NIPS link](#)

(already cited in last lectures)

- a good application of using FCFW when the linear oracle is expensive is given in this paper:  
Barrier Frank-Wolfe for Marginal Inference, R. Krishnan, S. Lacoste-Julien and D. Sontag, NIPS 2015  
<https://arxiv.org/abs/1511.02124>

- the generalization of the observation that Frank-Wolfe optimization sometimes reduced to the subgradient method on the primal is given in:
  - F. Bach. Duality between subgradient and conditional gradient methods.  
SIAM Journal of Optimization, 25(1):115-129, 2015  
[http://www.di.ens.fr/~fbach/sg\\_cg\\_fbach\\_siop.pdf](http://www.di.ens.fr/~fbach/sg_cg_fbach_siop.pdf)